

# Kleros

Long Paper v2.0.2

Clément Lesaege, William George, and Federico Ast

July 2021

## Abstract

Kleros is a decentralized decision protocol for use on smart contract platforms, which has been implemented on Ethereum. It acts as a decentralized third party capable of providing decisions on the correct result when applying a set of rules to questions ranging from simple to highly complex. This is achieved by using game theoretic incentives to have crowdsourced jurors analyze and rule on cases correctly. Hence, Kleros provides judgments in an inexpensive, reliable, typically fast, and decentralized way. Of particular relevance is the use of this protocol to dispute resolution, creating for a form of decentralized justice.

## 1 Introduction

*“Whoever controls the courts, controls the state”. Aristotle.*

The world is experiencing an accelerated pace of globalization and digitalization. An exponentially growing number of transactions are being conducted online across jurisdictional boundaries. If the blockchain promise comes to a reality, in a not so distant future, most goods, labor and capital will be allocated through decentralized global platforms. Disputes will certainly arise. Users of decentralized eBay will claim that sellers failed to send the goods as specified in the agreement, guests in decentralized Airbnb will claim that the rented house was not “as shown in the pictures” and backers in a crowdfunding platform will claim a refund as teams fail to deliver the promised results.

Smart contracts are smart enough to automatically execute as programmed, but not to render subjective judgments or to include elements from outside the blockchain. Existing dispute resolution technologies are too slow, too expensive and too unreliable for a decentralized global economy operating in real time. A fast, inexpensive, transparent, reliable and decentralized dispute resolution mechanism that renders ultimate judgments about the enforceability of smart contracts is a key institution for the blockchain era.

Kleros is a decision protocol for a multipurpose court system able to resolve every kind of dispute. It is an autonomous organization, implemented on Ethereum, that works as a decentralized third party to arbitrate<sup>1</sup> disputes in every kind of contract, from very simple to highly complex ones. Every step of the arbitration process (securing evidence, selecting jurors, etc.) is fully automated. Kleros does not rely on the honesty of a few individuals but on game-theoretical economic incentives.

---

<sup>1</sup>Here, a priori, the words “arbitrate” and “arbitration” are used in an informal sense. The status of Kleros decisions with respect to existing legal systems is a subject of research.

It is based on a fundamental insight from legal epistemology: a court is an epistemic engine, a tool for ferreting out the truth about events from a confusing array of clues. An agent (jury) follows a procedure where an input (evidence) is used to produce an output (decision) [42]. Kleros leverages the technologies of crowdsourcing, blockchain and game theory to develop a justice system that produces true decisions in a secure and inexpensive way<sup>2</sup>.

## 2 Previous Work: SchellingCoin Mechanism

Game theorist Thomas Schelling developed the concept of Focal Points, which are also called Schelling Points, [55] as a solution that people tend to use to coordinate their behaviour in the absence of communication<sup>3</sup>, because it seems natural or relevant to them. Schelling illustrated the concept with the following example: “Tomorrow you have to meet a stranger in NYC. Where and when do you meet him?”. While any place and time in the city could be a solution, the most common answer is “noon at the information booth at Grand Central Terminal”. There is nothing that makes noon at Grand Central Terminal a location with a higher payoff (any other place and time would be good, provided that both agents coordinate there), but its tradition as a meeting place makes it a natural focal point.

Based on the concept of Schelling Points, Ethereum founder Vitalik Buterin has proposed the creation of the SchellingCoin [15], a token that aligns telling the truth with economic incentives. If we wanted to know if it rained in Paris this morning, we could ask every owner of a SchellingCoin: “Has it rained in Paris this morning? Yes or No”. Each coin holder votes by secret ballot and after they have all voted, results are revealed. Parties who voted as the majority are rewarded with 10% of their coins. Parties who voted differently from the majority lose 10% of their coins.

Thomas Schelling [55] described “focal point(s) for each person’s expectation of what the other expects him to expect to be expected to do”. SchellingCoin uses this principle to provide incentives to a number of agents who do not know or trust each other to tell the truth. We expect agents to vote the true answer because they expect others to vote the true answer, because they expect others to vote for the true answer... In this simple case, the Schelling Point is honesty.

The majority votes \ You vote	YES	NO
YES	<b>+0.1</b>	<b>-0.1</b>
NO	<b>-0.1</b>	<b>+0.1</b>

Figure 1: Payoff table for a basic Schelling game

<sup>2</sup>For perspectives on criteria that can be used to evaluate whether Kleros-like systems based on such technologies can be considered to provide “decentralized justice”, see [7].

<sup>3</sup>Note that for applications of Schelling points in blockchain systems it is often impossible to guarantee that agents will not communicate as they tend to be pseudo-anonymous. While this is the case of Kleros, we will see that Kleros has an appeal system that incentivizes participants to agree with how potentially not-yet-determined agents in some future appeal round would decide, recovering a partial impossibility of communication. Furthermore, we are undertaking research on how to incentivize participants to not trust any communication they might have between each other, building off of work in [28] that argues that Schelling points also arise in such situations, see Section 4.9.

SchellingCoin mechanisms have been used for decentralized oracles and prediction markets [59] [49] [3]. We note academic work, developed concurrently to Kleros, that also attempts to apply these principles to dispute resolution<sup>4</sup> [21]. The fundamental insight is that voting coherently with others is a desirable behaviour that has to be incentivized. The incentives design underlying Kleros is based on a mechanism similar to SchellingCoin, slightly modified in order to answer to a number of specific challenges regarding scaling, subjectivity and privacy to make agents engage in adequate behaviour.

### 3 A Use Case

Alice is an entrepreneur based in France. She hires Bob, a programmer from Guatemala, on a P2P freelancing platform to build a new website for her company. After they agree on a price, terms and conditions, Bob gets to work. A couple of weeks later, he delivers the product. But Alice is not satisfied. She argues that the quality of Bob’s work is considerably lower than expected. Bob replies that he did exactly what was in the agreement. Alice is frustrated. She cannot hire a lawyer for a claim of just a couple hundred dollars with someone who is halfway around the world.

What if Alice and Bob had used Kleros Escrow and put a clause in the contract stating that, should a dispute arise, it would be solved by Kleros Court? After Bob stops answering her email, Alice taps a button that says “Send to Kleros” and fills in a simple form explaining her claim.

Thousands of miles away, in Nairobi, Chief is a software developer. In his “dead time” on the bus commuting to his job, he is checking the Kleros Court website (<https://court.kleros.io>) to find some dispute resolution work. He makes a couple thousand dollars a year on the side of his primary job by serving as a juror in software development disputes between freelancers and their clients. He usually rules cases in the Website Quality court. This court requires skills in html, javascript and web design to solve disputes between freelancers and their customers. Chief stakes 2000 PNK, the token used by Kleros to select jurors for disputes. The more tokens he stakes, the more likely is that he will be selected as juror. Assuming Chief rules well in the dispute, he gets his stake of 2000 PNK back afterwards in addition to a payment of arbitration fees, see Section 4.7.3.

About an hour later, an email hits Chief’s inbox: “You have been selected as a juror on a website quality dispute. Download the evidence [here](#). You have three days to submit your decision”. Similar email are received by Benito, a programmer from Cusco and Alexandru, from Romania, who had also staked their PNK in the Website Quality court. They were selected randomly from a pool of almost 3,000 candidates. They will never know each other, but they will collaborate to settle the dispute between Alice and Bob. On the bus back home, Chief analyzes the evidence and votes who is right.

Two days later, after the three juries have voted, Alice and Bob receive an email: “The jury has ruled for Alice. The website was not delivered in accordance to the terms and conditions agreed by the parties. A smart contract has transferred the money to Alice”. Jurors are rewarded for their work and the case is closed.

---

<sup>4</sup>Compared to Kleros, the proposal of [21] uses a somewhat different structure for how arbitrators stake to be selected for, vote on, and then receive fees for disputes, with the Schelling game occurring in an additional “validation” phase. Other notable differences are that, while the proposal of [21] does not include a mechanism similar to the court tree we describe in Section 4.3 for global 51% attack resistance, it does include a novel proposal for a “forum” in which community members are incentivized to participate beyond the direct providing of votes in disputes, for example by submitting proposals for changes in governance or new template contracts.

The image displays a user interface for a dispute resolution process, divided into two main sections: an evidence review area and a voting area.

**Evidence Review Area:**

- Header:** "Evidence" with a folder icon and a scroll-to-bottom button.
- Timeline:** A vertical line with four key events: "Latest", "Jurors ruled in favor of the challenger", "Challenged", and "Dispute Created".
- Evidence 1:** A card titled "Here you can find the documents." with the text "The documents can be seen attached to support the argument in favor of Bob." and a submission footer: "Submitted by 0x1234...4567 09 Jan 2019 05:00:37 GMT".
- Evidence 2:** A card titled "It does not comply. It seems fake." with the text "This evidence shows that the previous document is fake according to the original agreement." and the same submission footer.
- Evidence 3:** A card titled "Here you can find the documents." with the text "The documents can be seen attached to support the argument." and the same submission footer.
- Evidence 4:** A card titled "Here you can find the documents." with the text "The documents can be seen attached to support the argument in favor of Alice." and the same submission footer.
- Footer:** "Start" and "Token Submission" buttons, along with a scroll-to-top button.

**Voting Area:**

- Header:** "Dispute History" with a scale icon and a scroll-to-top button.
- Section Header:** "What is your decision?" with a scroll-to-top button.
- Text Input:** A large text area with the placeholder "Justify your vote here ...".
- Buttons:** Three buttons are present: "Vote for Alice" (orange), "Vote for Bob" (orange), and "Refuse to Arbitrate" (grey) with an information icon.

Figure 2: The evidence before a juror as she makes her vote.

## 4 Kleros Mechanism Design

In this section, we detail the architecture of Kleros<sup>5</sup>.

### 4.1 Arbitrable and Arbitrator Contracts

Kleros is an opt-in court system. “Arbitrated” or “arbitrable” smart contracts have to designate Kleros as their arbitrator. When they opt-in, contracts creators choose how many jurors and which court will rule their contract in case a dispute occurs, see Section 4.3.

The Kleros team has developed a number of standard contracts using Kleros as a dispute resolution mechanism. Moreover, we have proposed standards [43] [61] that would allow others to develop other contracts in a way that does not require anticipating which dispute resolution mechanism will ultimately be used. This standardization allows parties that regularly require dispute resolution services to easily switching between dispute resolution providers.

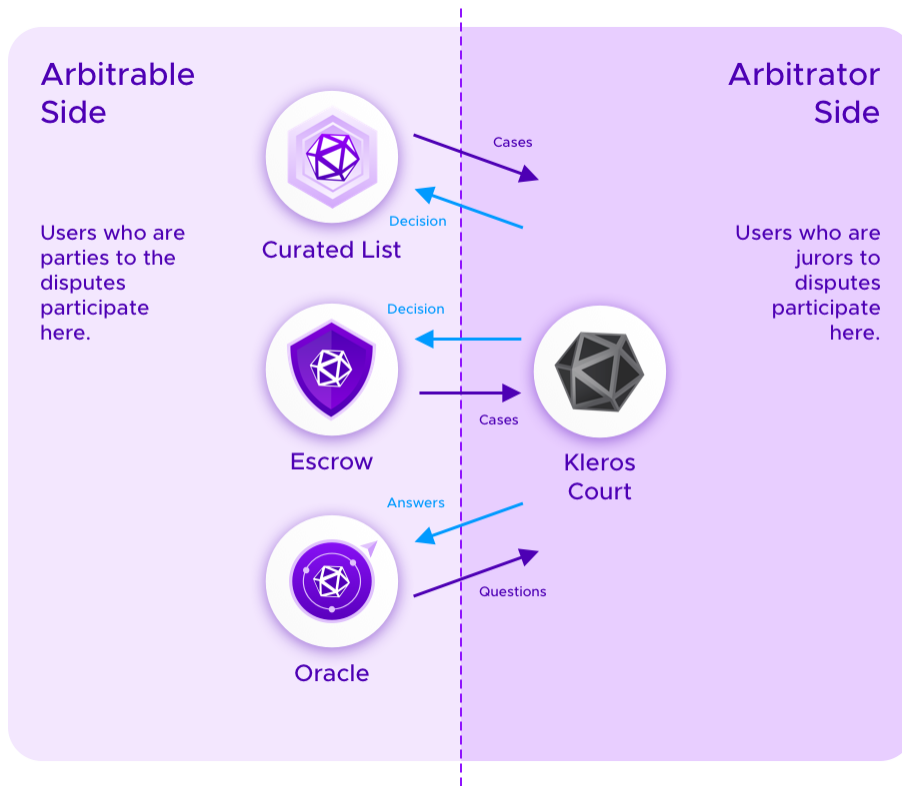


Figure 3: “Arbitrable” contracts, namely contracts that might require dispute resolution, are represented on the left. These include applications such as oracles, where users have some dispute over the true state of some off-chain information; curated lists, where disputes arise between users who argue that an entry satisfies some fixed set of criteria that qualify it to be on a list while other users disagree; and more general escrow applications, where a smart contract holds something of value pending the resolution of a dispute. See Section 6 for more information on these examples. “Arbitrable” contracts designate an “arbitrator” contract such as that of the Kleros Court, represented on the right, to provide rulings on any eventual disputes.

<sup>5</sup>Throughout this work we will indicate both the current state of Kleros as it is implemented at the time of writing as well as proposed mechanisms that we have researched for future versions. As we continue to research these issues, the proposed future mechanisms described below are subject to change.

Contracts can allow for a variety of possible behaviour in cases that do not require arbitration, particularly when there is unanimous consent of parties. For example, the current Kleros Escrow includes a limited settlement system.

Moreover, contracts will specify the options available for jurors to vote. In the introductory example, options may be: “Reimburse Alice”, “Give Bob one extra week to finish the website” and “Pay Bob”. The smart contract will also specify the behaviour of the contract after the ruling is done for each possible option. In the example:

- “Reimburse Alice” transfers funds to Alice’s address.
- “Give Bob one extra week to finish the website” advances the timers for how long Bob has to finish the website one week, i.e. blocks Alice from creating new disputes during this time. Furthermore, the smart contract might be written so that if this option has been chosen once, it cannot be selected in further disputes.
- “Pay Bob” transfers funds to Bob’s address.

In general, any finite list of options can be proposed to jurors<sup>6</sup>. Sets of options having a different structure can be accommodated in some circumstances to the degree that they can be “discretized”. For example, we have researched mechanisms by which jurors can choose a real-number value [31], as one might want for cases where one party receives a percentage refund. For examples, jurors might rule to refund 75% of an amount at stake to Alice while paying 25% to Bob.

## 4.2 Underlying Blockchain(s) and Scaling Solutions

Kleros is currently implemented on Ethereum; however, this protocol could be implemented on any blockchain with adequately expressive smart contracts. For example, an xDai version of Kleros is in the process of being launched [50]. Note, in the future, one could imagine a mechanism where cases from such sidechain courts could be appealed to the Ethereum version of Kleros. Indeed, one wants the security that comes with being able to appeal to the most robust version of the court in the event of an attack or a particularly contentious case; see Section 4.8 for more information on the appeal mechanisms.

Moreover, it is possible to have arbitrable and arbitrator contracts that are not in the same consensus, namely for these smart contracts to be on different chains or, for example, for one contract to be on an optimistic rollup on Ethereum, while the other contract is on L1 or on a different optimistic rollup [40]. This raises questions regarding communication between the arbitrable and arbitrator contract so that disputes can be raised and rulings can be enforced in an efficient and timely manner. Indeed, while contracts that are deployed on the same consensus can instantly interact atomically, this is not the case when contracts are on different chains or on different rollups. These questions are particularly relevant for Kleros, as it is a protocol that is designed to interact with a variety of different dapps that require dispute resolution, which may be on deployed on a variety of platforms.

The communication mechanisms available between different chains and different rollups vary significantly. In some cases, particularly when either the arbitrable or arbitrator contract is on an optimistic rollup, and also in the case of some cross-chain bridges, there exist “slow” means of communication

---

<sup>6</sup>One option that is always presented to jurors is to “refuse to arbitrate”. For the sake of incentivization of jurors, see Section 4.7.3, this option is treated like any other option and hence jurors have an economic incentive to vote for it if and only if they believe that it will be the winning option. A court policy of the Kleros General Court, and hence all Kleros courts, see Section 4.3, is that jurors should vote “refuse to arbitrate” if the decision is being used as part of illegal activities.

that require a delay of several days, corresponding to challenge periods where third parties can flag malicious communication, before messages sent from one contract can be processed by the other [40]. See Figure 4 for a mechanism that Kleros can use to facilitate faster communication between arbitrable and arbitrator contracts in these cases, where users make claims about the state of the other contract. These claims are challengeable and deposit-backed, similar to state updates transactions in optimistic rollups [40], however their challenge periods can be more adapted to the needs of arbitrable applications, so that one must only wait for the full period required by the “slow” communication mechanism in cases of conflict.

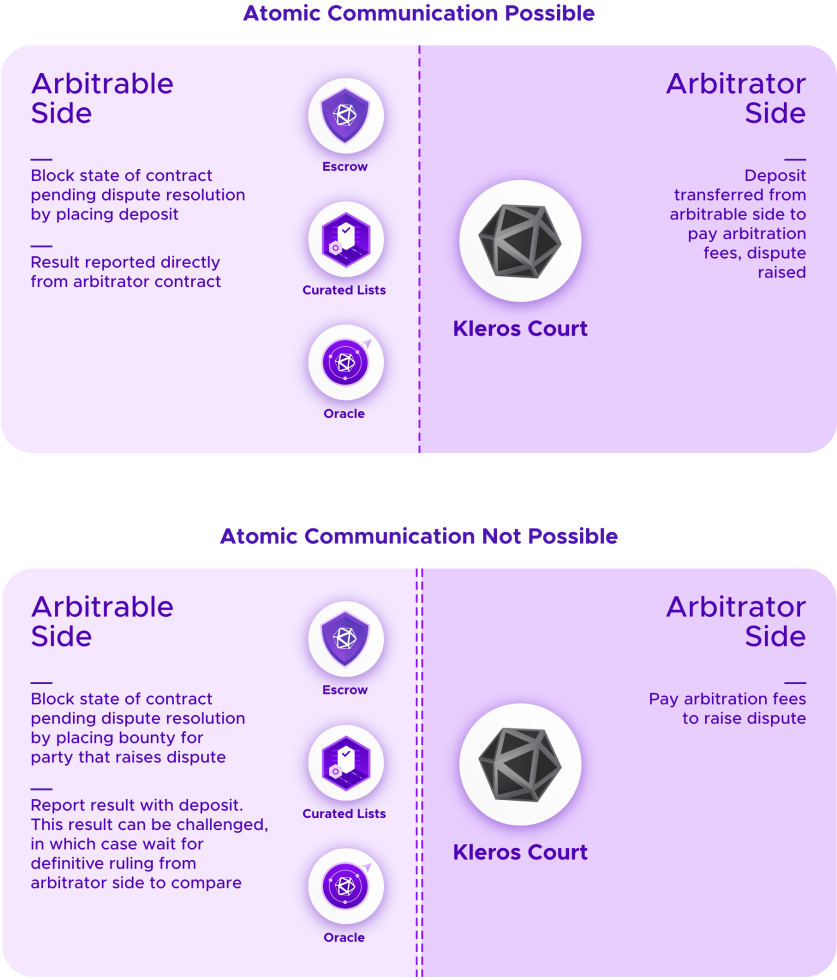


Figure 4: When atomic communication between arbitrable and arbitrator contracts is possible, a single transaction can both pause the state of the arbitrable application and pay arbitration fees to raise a dispute on the arbitrator application. Similarly, results of a dispute can be reported to the arbitrable application instantly. When atomic communication is not possible, but there is nevertheless a “slow” means of communication, such as is available when bringing information from an optimistic rollup to L1, one can have a schema where one blocks the state of the arbitrable application by placing a deposit that is somewhat larger than the required arbitration fees. Then this deposit serves as a bounty to parties who would pay the arbitration fees to raise a dispute on the arbitrator side. Moreover, one can allow users to submit challengeable information to the arbitrable contract along with a deposit, allowing for communication that “usually” requires only intermediate length delays and only defaults to the “slow” communication channel in case of a challenge.

Note, that schema for non-atomic communication of Figure 4 requires extra deposits to be made that are not necessary in cases where the arbitrable and arbitrator contract can directly communicate atomically. Nevertheless, an observer following both chains/both rollups can make these deposits in a risk-free fashion. Hence, by providing an additional fee, liquidity providers can be incentivized to provide this service<sup>7</sup> <sup>8</sup>.

### 4.3 Court Tree

When creating an arbitrable contract, parties should choose a type of court specialized in the topic of the contract. A software development contract will choose a software development court, an insurance contract will select an insurance court, etc.

In parallel, when registering as jurors, users start in the General Court and follow a path to a specific court according to their skills<sup>9</sup>. Each token holder can register a given token in at most one child court of each court in which they have staked tokens<sup>10</sup>. Figure 5 shows the current court structure with examples of allowable registrations.

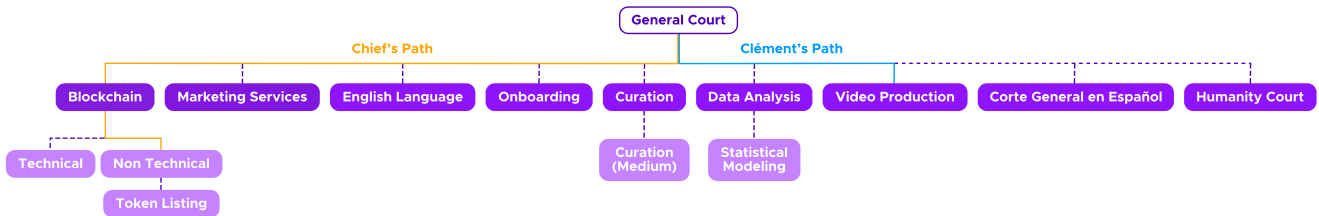


Figure 5: The current court tree from which smart contract creators must select a court. New courts will be added as Kleros is adopted to resolve additional types of disputes. The mechanism to add new courts is described in Section 4.12. Example of paths chosen by jurors in the court system are shown. Clément can be drawn as juror in the General Court and in the Video Production Court. Chief can be drawn as juror in the General Court, in the Blockchain Court, and in the Blockchain Non-Technical Court.

Asking jurors to choose between courts incentivizes them to choose the courts they are the most skilled at. If they were able to choose every court, some would choose all of them to get the maximum amount of arbitration fees from their tokens.

Each court has some specific features regarding policies. Also a number of parameters are chosen on a court by court basis including session time, cost, number of drawn jury members and tokens staked. We will consider these system parameters in detail below, see Section 4.7.3.

<sup>7</sup>Moreover, in some cases this structure can be simplified. For example, if one of the contracts is on Ethereum L1, communication from that contract to a contract on an optimistic rollup is more facilitated than the inverse. Additionally, when the dispute is structured in such a way that a party would be incentivized to raise a dispute without needing the reward of a bounty, for example because this party has an interest in the dispute and would lose by default if the dispute if arbitration fees are not paid, then paying deposits on the scale of the arbitration fees to both the arbitrable contracts and the arbitrator contract can be avoided.

<sup>8</sup>The next version of Kleros is intended to be implemented on a Layer 2 solution on Ethereum using optimistic rollups. However, we are following developments in zk-rollup technology closely, and this may be used to the degree that zk-rollups for general purpose smart contracts are available.

<sup>9</sup>In the language of graph theory, the structure of the set of courts forms an arborescence with the General Court as the root.

<sup>10</sup>Moreover, token holder *must* be staked in the parent court to stake in the child court, so when staking in a given court a token holder is staked in all the courts on the path from it to the General Court.



## 4.4 Raising a Dispute

The nature of an arbitrable contract determines its default behaviour and conditions under which a dispute can be raised<sup>11</sup> Typically, parties to an agreement can raise a dispute if there is a disagreement, or accept a default behaviour if there is agreement.

In the event of a dispute, parties can provide evidence and arguments on behalf of their case during an evidence period. This evidence conforms to the ERC 1497 standard [61], which sets forth requirements for how this evidence is organized and how it triggers smart contract events, providing for interoperability across arbitrator applications.

## 4.5 Drawing Jurors

### 4.5.1 Staking

Users have an economic interest in serving as jurors in Kleros: collecting arbitration fees in exchange for their work. Candidates self-select to serve as jurors by staking a Kleros crypto-token, called PNK<sup>12</sup>. The probability of being drawn as a juror for a specific dispute is proportional to the amount of tokens a juror stakes. The higher the amount of tokens she stakes, the higher the probability that she will be drawn as juror. Staking PNK signals availability to be drawn as a juror; users that do not stake PNK do not have the chance of being drawn. This prevents inactive jurors from being selected.

PNK plays three key functions in the design of Kleros.

- First, it protects the system against Sybil attacks [25], namely attacks where malicious parties create many addresses in order to be drawn a high number of times in a given case. Depending on the mode of the juror selection process, see Section 4.5.2, this protection can be direct, where such an attacker would be required to stake as many tokens as those staked by honest candidate jurors in order to have a high likelihood of success, see Section 4.5.2.1. Alternatively, this protection can come indirectly via Sybil resistant mechanisms whose own security follows from that of Kleros courts using the preceding mode, see Section 4.5.2.2.
- Second, PNK provides jurors the incentive to vote honestly by making incoherent jurors, i.e. jurors whose votes do not agree with the ultimate ruling, pay part of their stake to coherent ones, see Section 4.7.3.
- Finally, PNK can be “forked” in a way that creates parallel versions of Kleros, serving as a fallback defense in the event of successful 51% attacks, see Section 4.10.

### 4.5.2 Jury Selection

Once candidate jurors have expressed their availability to be drawn as jurors by self-selecting into specific court by staking their tokens in that court, the final selection of jurors is done randomly<sup>13</sup>. For details on the process used to generate random numbers for this selection, see Section 4.5.3. The probability of being drawn as a juror is proportional to the amount of staked tokens.

Theoretically, drawing a candidate purely in proportion to her staked tokens allows for the possibility that the same candidate may be drawn more than once for a specific dispute. However, in practice

---

<sup>11</sup>Such as time limits, or fees required to Kleros jurors must be paid, see Section 4.6.

<sup>12</sup>The ticker of this token, PNK, is a reference to the pinakion, the bronze plaque that each Athenian citizen used as an ID. The pinakion was used as a token for jury selection in Athens popular trials.

<sup>13</sup>Note that random draw has a long history of use in public decision making. In addition to its use in jury selection in many countries today, it was, for example, used in the selection of public office-holders in ancient Athens and Renaissance Venice [22]. For reflections on the legitimacy of such processes, see [26] [58]

this is unlikely for disputes with a small number of jurors in popular courts with diverse staking pools. Nonetheless, we consider two modes for drawing jurors according to how Sybil resistance is achieved and the consequences this has for whether addresses are allowed to be selected multiple times for the same dispute.

Which mode is used by a given subcourt is controlled by a parameter that can be adjusted through governance, see Section 4.12.

### 4.5.2.1 Sybil Resistance through Tokens

Note that, absent alternative Sybil resistant mechanisms such as that described in Section 4.5.2.2, restricting the number of votes per address would allow for malicious parties to split their token holdings over multiple addresses and gain an advantage compared to honest parties. To eliminate these opportunities for manipulation, according to this mode, the possibility to draw the same candidate two (or more) times is allowed. The number of times a user is drawn for a dispute (called its weight) determines the number of votes she will get in the dispute and the amount of tokens she will win or lose during the token redistribution. Thus, an attacker gains no particular advantage by dividing her tokens over multiple addresses compared to participants who hold all of their tokens on a single address. As a result, in order for an attacker to have a high likelihood of receiving more than half of the votes in a given case under this mode, she must have actually have a majority of the tokens staked, see our discussion of 51% attacks in Section 4.11.1.



Figure 6: Imagine that 6 token owners staked 10,000 in total with the above distribution. Then 5 tokens are drawn: numbers 2519, 4953, 2264, 3342 and 9531. As a result, token owners B, C and F are drawn with a weight of 1. Token owner D is drawn with a weight of 2.

See Figure 6 for an example of juror selection. Note that staked PNK (except for a part of those staked by incoherent jurors) can be taken back after the court reaches a final decision.

### 4.5.2.2 Sybil Resistance Using Proof of Humanity

In future versions of Kleros, an additional mode of Sybil resistance will be available. Here, in order to be selected as a juror in a court that uses this mode, the juror’s address must be registered on Proof

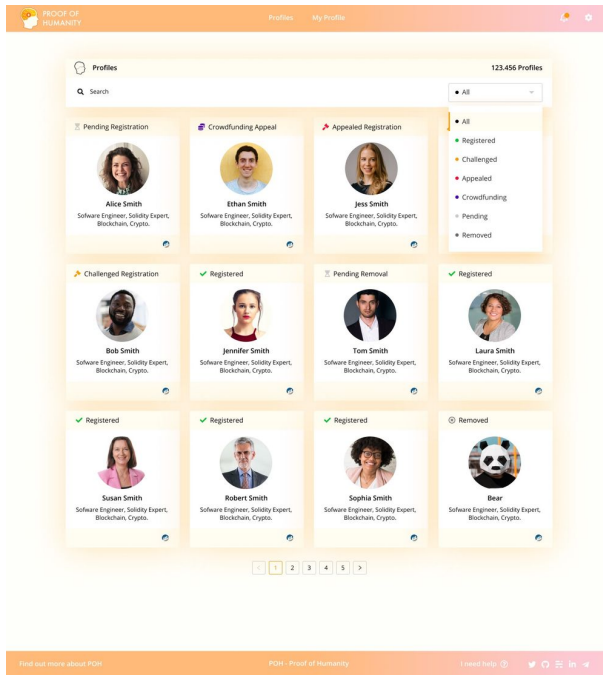


Figure 7: Proof of Humanity is a curated list of distinct humans, where disputes regarding whether a submission is a duplicate or an invented person are resolved by Kleros. In future versions of Kleros, certain courts will follow a mode where in order to become a juror in that court, a user must be registered on Proof of Humanity, and then the Sybil resistant properties of this list prevent malicious users from being able to obtain a large number of votes through manipulation.

of Humanity [38], a Sybil resistant list of humans, see Figure 7. This list is curated via a mechanism that itself uses Kleros courts, specifically the “Humanity Court” and its appeal court(s), to rule on whether entries are duplicates or otherwise malicious.

In exchange for the additional UX friction that requiring this registration creates, this model allows one to limit each address to at most a single vote while nonetheless eliminating the possibility for malicious parties to gain an advantage against honest users by splitting their token holdings over multiple addresses.

The mode of Section 4.5.2.1 is particularly useful for courts that resolve disputes regarding Proof of Humanity itself, such the Humanity Court and its appeal court(s), notably the General Court. Indeed, using Proof of Humanity to provide for Sybil resistance in cases regarding itself could lead to circular security guarantees and limit the ability of Proof of Humanity, and by extension the courts that use it for Sybil resistance, to recover from attacks. Moreover, this mode would also be natural for courts where an adequate pool of qualified jurors has staked, but is nonetheless not registered in Proof of Humanity.

### 4.5.3 Random Number Generation

In order to draw jurors, a random number generation process that is resistant to manipulation is required. Using a protocol for creating a random number between two parties [8] does not work. An attacker could create disputes between herself, select herself as a juror multiple times and select another victim juror. She would then coordinate her own votes in a way that they would be considered coherent but not those of the victim in order to steal tokens from the victim when PNK are redistributed, see the discussion of the incentive system in Section 4.7.3.

### 4.5.3.1 Random Number Generator Using Blockhashs

Currently, the random numbers used to select jurors are drawn from blockhashes of Ethereum blocks. While these values remain impossible to predict in advance, miners can opt not to release a block that would result in random numbers unfavorable to them at the expense of forfeiting a block reward.

### 4.5.3.2 Random Number Generation Using Threshold Signatures

Another possible random number generator is one based on threshold signatures. Here the idea is that a number of “key-keepers” generate shares of a  $t$  of  $n$  threshold signature key, with a corresponding public verification key. For example, one can use the threshold version [10] of the BLS signature scheme [11]. Then, in order to generate a random value, a given seed is initialized, and the key-keepers collectively produce the corresponding signature on this seed, which can be verified as being correct using the public verification key. Then as long as the threshold of  $t$  key-keepers participate, the resulting signature will be a fixed value, regardless of the actions taken by the  $n - t$  remaining key-keepers. Namely, no group of less than  $n - t + 1$  key-keepers can change the resulting random value nor prevent this process from returning a value. Moreover, no group of less than  $t$  colluding key-keepers should be able to recover the random value prematurely. See [24] for work in this direction; particularly the Chainlink Verifiable Random Function (VRF) [23] uses this approach.

### 4.5.3.3 Random Number Generation Using Sequential Proof-Of-Work

Another approach that may be available in the future would be to use a Verifiable Delay Function [9]. One approach to produce such a function is based on sequential Proof-Of-Work, such as by using a scheme similar to Bünz et al. [20]. The idea of this scheme is that a minimum amount of time, based on the fastest available hardware to make certain non-parallelizable computations, should be required to compute the random value. Then no actor should be able to predict the generated random value faster than this lower bound.

We briefly summarize such a scheme:

1. **Initialization:** One starts with  $\text{seed} = \text{blockhash}$  and let all parties input a value  $\text{localRandom}$  to change the seed such that  $\text{seed} = \text{hash}(\text{seed}, \text{localRandom})$ . This allows any party to change the seed. It is important that the seed is not chosen by any one party. Using the above process, every party can change the seed, but not choose it, because choosing a particular  $\text{seedAttack}$  would require the attacker to determine  $\text{localRandom}$  such that  $\text{hash}(\text{seed}, \text{localRandom}) = \text{seedAttack}$  which is difficult due to the preimage resistance of cryptographic hash functions<sup>14</sup>.
2. **Computing the master random value:** Every party who has a stake in the random number runs sequential Proof-Of-Work on the seed. Starting with  $h_0 = \text{seed}$ , they compute  $h_{n+1} = \text{hash}(h_n)$  up to  $h_d$  where  $d$  is the difficulty parameter. Computing  $h_d$  takes time and assures that a certain amount of time passed between someone gets the knowledge of the seed and that she gets the result. The difficulty  $d$  is fixed such that no hardware can compute  $h_d$  during the time of the initialization phase. Because one needs the result of the previous step before starting the next one, this process cannot be parallelized. This means that no party will be able to obtain the results significantly faster than the others.

---

<sup>14</sup>Note that the protocol that we present here is adapted to work for Proof-Of-Stake blockchains as well as for Proof-Of-Work blockchains. In Proof-Of-Work blockchains, as the blockhash remains impossible to exactly predict, one can remove the initialization step and only use the blockhash as a seed. However, Ethereum has planned to switch to Proof-Of-Stake.

3. **Getting the results on the blockchain:** Every party can post the  $h_d$  with a deposit they found. Then other parties can disprove results which are wrong using interactive verification [51]. It consists of a dichotomic search on the results of the attacker. If an attacker submits a false  $h_d$ , an honest party can ask her for her  $h_{d/2}$  value. If she gives the wrong value, there is an error in the attacker values between  $h_0$  and  $h_{d/2}$ . If she gives the right value, there is an error between  $h_{d/2}$  and  $h_d$ . Either way, the search space is divided by two. The honest party continues this process on a reduced space (where the error is) until two values are left. Then the honest party can exhibit  $x$  such that  $h_{x+1} \neq \text{hash}(h_x)$  in the attacker answer which invalidates her answer. Parties whose answer is invalidated lose their deposit. Part of it is burnt and the other part is given to the party that invalidated them. Note that the number of interactions required to invalidate a false result is only  $O(\log(d))$ .
4. **Getting all random values:** After the honest parties have invalidated the results, there is only the correct result  $h_d$  left. From this master random value one derives all the random values such that  $r_n = \text{hash}(h_d, n)$ .

The output of this process is a random number as long as there is at least one honest party. Computing the sequential Proof-Of-Work and the interactive verification takes time. But for most disputes waiting a few hours from the moment the dispute starts and the moment jurors are drawn will not be a problem. However, for some courts with a particularly low session time (for example, a court solving disputes related to content moderation on a decentralized social media platform) this random number generation method could be too slow.

Which random number generation technique is used will depend on what methods are available on a given underlying blockchain platform, see Section 4.2. For example, blockhashs are not a notion native to most optimistic rollup platforms [40] and are hence not available for random number generation. On the other hand, threshold signature based schemes such as that of the Chainlink VRF [23] are planned for such platforms [47]. Verifiable delay functions, such as that described in the above sequential Proof-Of-Work scheme, may provide an appealing long-term solution.

## 4.6 Arbitration Fees

Kleros uses arbitration fees in order to compensate jurors for their work. These fees also make it more difficult for an attacker to spam the system by creating frivolous disputes and/or appealing as these actions require paying arbitration fees. Each juror who is coherent with the final ruling will be paid a fee determined by the court where the dispute is solved. The arbitrable smart contract will determine which party will pay the juror fee; this can vary from one application to another.

The rules can be simple. For example, they may require the party creating the dispute to pay the fee. However, we may think of more complex rules to create better incentives. For example, one could require each party will deposit an amount equal to the juror fee in the smart contract. If one party fails to do so, the smart contract will consider that the court ruled in favor of the party who deposited the juror fee (without even creating a dispute in the court). If both parties deposit the funds, the winning party will be reimbursed when the dispute is over.

## 4.7 Voting and Incentivization

### 4.7.1 Voting Process

Jurors assess evidence that has been submitted, typically by the parties to the dispute, and are provided with court policies, comparable to juror instructions<sup>15</sup>, on how they should reason based on that evidence.

Then each juror commits [13] to a vote that reflects her ruling on the case. Namely, she submits  $\text{hash}(\text{vote}, \text{salt}, \text{address})$ <sup>16</sup>. In future versions of Kleros, the juror will be able to resubmit a commitment during the voting period, with a new salt, either to reconfirm an existing vote or to change a vote. When the vote is over, jurors reveal  $\{\text{vote}, \text{salt}\}$ , and a Kleros smart contract verifies that it matches the (final) commitment received for her vote. Jurors failing to reveal their vote are penalized, see Section 4.7.3. For more information about this process of committing to a vote via a hash value, particularly in the context of limiting the circulation of information on the votes of other jurors during the voting period, see Section 4.9.

### 4.7.2 Vote Aggregation

After the end of the voting period, votes are revealed by jurors. Jurors that fail to reveal their vote are penalized. Finally, votes are aggregated according to a predetermined voting rule resulting in an option that is considered the winner.

#### 4.7.2.1 Current Voting System

When jurors are presented with a binary choice, it is natural to use the Plurality, or “first-past-the-post”, voting system<sup>17</sup>. Currently, Kleros employs the Plurality system even when there are more than two choices. (Specifically, the winning choice is that with a plurality of the votes in the last appeal, see Section 4.8).

When there are more than two options, under a Plurality voting system, the following may occur:

- If there are many very similar honest options (or “clones”), they will divide the votes of jurors that are attempting to vote honestly, decreasing the probability that any one of them wins. Anticipating this effect, jurors might instead vote for a distinguished but dishonest choice. For example, imagine that in our contractor use-case above there were several different options that gave Bob another week, another eight days, or another nine days respectively. Then the collective odds of these options being chosen could fall below the odds of a single “give Bob more time” option, see Figure 8.
- To the degree that no single option is likely to receive more than 50% of the votes, this lowers the bar for the number of votes that attackers need to corrupt to pass a dishonest result, see Figure 9.

Considering these issues, one might expect Plurality voting to still produce generally “honest” results when one single choice has a very clear, winning case, which essentially binarizes the choice. However, Kleros should be able to cope with nuanced cases involving many choices.

---

<sup>15</sup>These policies vary by court, see Section 4.3 and can be changed by the governance procedure, see Section 4.12.

<sup>16</sup>Throughout this paper we use hash referring to a cryptographic hash function, in Ethereum the one used is `keccak256`.

<sup>17</sup>Plurality, or “first-past-the-post”, is the voting system that in which voters express a vote for only one candidate, and then the candidate that receives the largest number of votes is selected, even if this candidate does not receive a majority of the total votes due to there being more than two candidates.

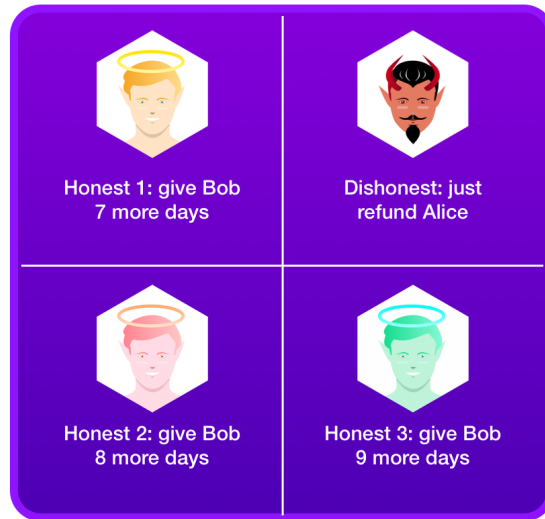


Figure 8: In the Plurality voting system, jurors can only vote for one option. Particularly, a voter cannot cast a vote such as “take one of the more time options, it doesn’t matter which”. Then if jurors are presented with a collection of similar, honest choices along with a single dishonest choice, the dishonest choice may seem distinguished and become the Schelling point.

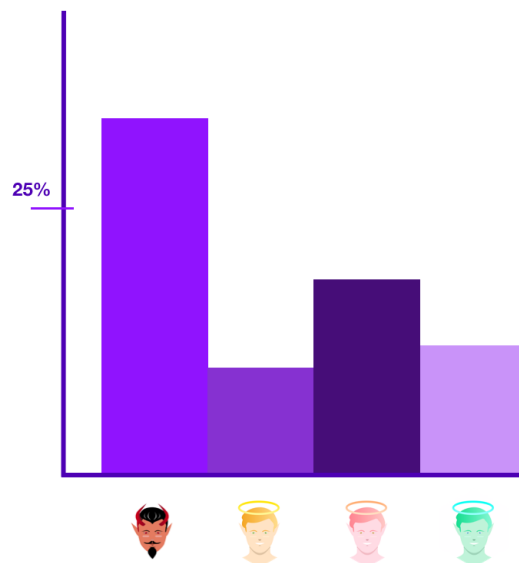


Figure 9: If the vote is split between several “honest” options, the attacker does not need to corrupt 50% of the vote in a Plurality system to have a malicious option adopted.

### 4.7.2.2 Social Choice Theory and Future Voting System(s)

In this section, we consider a number of desirable properties for a vote aggregation rule in a system such as Kleros, remarking which rules satisfy or do not satisfy these properties. These vote aggregation rules assume that jurors submit more nuanced information about their opinions on a case than just a single outcome that they consider to be the best; such votes often take the form of a ranked list of outcomes  $a_1 \geq a_2 \geq \dots \geq a_n$ . In Section 4.7.3 we will discuss possible payoff structures that can be used to incentivize participants in such votes. Note that for some of the properties we consider, an analysis of how a given vote aggregation rule performs will depend on the choice of incentive system. Indeed, while we discuss vote aggregation and payoff structures separately for clarity, in fact these choices are deeply interwoven and tradeoffs regarding them should be considered together. Also note that, while some of the properties we consider have standard definitions, others will be somewhat subjective.

In future versions of Kleros, the smart contracts making up the court will be more modular. In particular, it will be possible for smart contracts with different voting and incentive mechanisms to be written that then can be executed by an existing court contract. Which voting and incentive systems are allowed will be determined by the governance mechanism, see Section 4.12. For different applications it may make sense to prioritize different properties when considering the tradeoffs that we present.

Ideally, the vote aggregation system should have the following properties:

- Clone independent - From the point of view of a voting system, this is a standard property considered in social choice theory that informally means that having a “clone set” of multiple similar options does not increase or decrease the winning chances for other options outside of the clone set, see for example [60] for a formal definition. Using a clone independent aggregation rule allows parties to produce arbitrable contracts without having to worry that presenting two or more similar options to the jurors might lead to vote splitting in or against their interests.

Examples of clone independent voting systems are Instant-Runoff (IRV), Ranked Pairs, and Schulze. The results in the clone independence column of the table below, in the standard sense as a voting system, are all presented in [60] and [56]. Note, however, that in the situation we are considering, a more robust notion of clone independence would demand that not only the voting system be clone independent, but that also the addition of a clone not affect the expected payoffs that voters receive, or at least not affect which votes provide the optimal payoffs, up to reordering of alternatives within a clone set. Future research may formalize this criterion.

Nonetheless, we have analyzed these systems heuristically/numerically and we have noted that IRV-type systems, while being clone independent as voting systems, introduce a bias in the rewards of clones for the types of payoff systems we will consider in Section 4.7.3. Indeed, note that if a voter highly ranks the clone set, then under an IRV-type system which clone the voter ranks first has more of an influence on which clone wins than if the voter ranks the clone set lower. Consequently, incentive systems that reward or penalize voters based on where they place the ultimate winner relative to other choices allow voters that rank a set of clones highly to capture a slightly larger percentage of rewards after the addition of an option to the clone set. We have, at least heuristically, not observed such a bias for voting systems such as Ranked Pairs and Schulze, where the aggregation rule is based on the graph of the relative strengths of pairwise duels.

- Satisfy the Condorcet criterion - A voting rule is said to be a Condorcet method if whenever there is an option  $w$  such that more voters rank  $w$  higher than  $a$  for all other options  $a$ , then  $w$  wins, see [12]. Note, however that such a “Condorcet winner” does not necessarily exist for all



sets of votes expressed by jurors. The idea that, if there is an option that wins “head-to-head” against all others, then that option should be selected, is fairly intuitive. Hence, if Kleros cases often have Condorcet winners, then they are particularly straightforward to mentally simulate for jurors. Furthermore, the Condorcet criterion, namely to prefer options that have a consensus of the population against each other option, corresponds to certain notions of “fairness”.

Moreover, we see that the Condorcet property has positive effects on attack resistance when combined with the types of incentive systems that we consider in Section 4.7.3.2. Indeed, an attack that attempts (and fails) to dislodge a Condorcet winner by reversing a duel between that alternative and another is something that is well captured, and hence penalized, by these incentive systems. Ranked Pairs and Schulze are examples of Condorcet methods. WoodSIRV is essentially a version of IRV that is Condorcet-ized by checking at each voting round whether there is a Condorcet winner and selecting it if available.

- Resistant to attacks - Note that attack resistance is essentially economic. For example, one wants the number of votes that would need to be changed to result in a “dishonest” outcome winning to be high, with the rationale that this increases the cost of attacks, such as bribes that would be necessary to change those votes. Hence, attack resistance touches on certain standard voting systems properties that have been studied in social choice theory, such as the Later No Help and the Participation criteria [12], in combination with analysis of the penalties and rewards laid out in the incentive system, see Section 4.7.3.

While one can exhibit attacks on individual systems, it is challenging to produce proofs that no attacks exist<sup>18</sup> The discussion in Section 4.7.2.1 already shows that the bar for a 51% attack is lowered on the Plurality system in cases where no single (honest) option receives more than 50% of the votes. In the Borda system, if the “honest” answer is  $a$  and an attacker wants the option  $b$  to win, she can submit many votes where  $a$  is ranked first and  $b$  is ranked second, and a small number of votes where  $b$  is ranked first and  $a$  is ranked last. Using the incentive system(s) that will be described in Section 4.7.3, this attack has limited cost and risk for the attacker whereas for appropriate choices of parameters, Instant-Runoff and Ranked Pairs seem to be more resistant<sup>19</sup>.

It is also worth keeping in mind which types of manipulative behaviours the incentive systems we consider in Section 4.7.3.2 are capable of detecting and penalizing so as to avoid choosing an aggregation function that is vulnerable to manipulation that the incentive system does not penalize. We have already noted above that an attack that prevents a would-be Condorcet winner from winning is well detected by the incentives we consider. On the other hand, according to Arrow’s Impossibility Theorem [5], for all non-trivial voting systems, there will be situations where the rankings of “irrelevant alternatives” will affect the results. Manipulations involving such irrelevant alternatives are not well captured by the incentive systems we consider in Section 4.7.3<sup>20</sup>. However, note that as Instant-Runoff satisfies Later No Harm and Later No Help, such manipulations are only relevant to the degree that they involve alternatives that are already

---

<sup>18</sup>Hence the question marks on the claims of better attack resistance for Instant-Runoff, WoodSIRV, Ranked Pairs, and Schulze in the table below.

<sup>19</sup>All of our claims on attack resistance should be viewed as based on using the incentive systems of Section 4.7.3. It is possible with another incentive system our conclusions would be different. Also note the role of  $\beta$  discussed in Section 4.7.3; for  $\beta = 0$  in weighting systems 2 and 3, the economic cost in lost deposits of the attack described on Borda is comparable to that of an attack on Instant-Runoff where the attacker bribes many voters to place  $b$  first and  $a$  second, though the attack on Instant-Runoff requires convincing a large number of jurors to accept a small in-protocol penalty, and the attack on Borda involves convincing a small number of jurors to accept a large penalty. However, with other choices of  $\beta$ , this attack on Instant-Runoff becomes more expensive while the attack on Borda does not.

<sup>20</sup>This is a consequence of these systems having the Winner First, Maximum Payout property that we discuss in

ranked higher than the alternative that would otherwise win. Then, in this system, one would expect that in order to produce a dishonest answer, the attacker would need to corrupt many votes willing to place a dishonest answer over the “honest” answer.

Thus, we argue that voting systems such as WoodSIRV, which are Condorcet, but resolve Condorcet paradoxes with an IRV like step, seem to offer the better attack resistance as they combine these layers of defense, requiring an attacker to bear the cost of preventing a Condorcet winner if one exists, and then to the degree that a winner is manipulated, when finds oneself in a Condorcet paradox, only manipulations by voters who already receive a significant penalty are useful, limiting effective attacks to something heuristically resembling the coalition required for a 51% attack<sup>21</sup>

- Not require too much code complexity nor gas - Even on a scaled version of Ethereum, gas costs would likely remain a usability concern. Moreover, in any smart contract platform, reducing code complexity is beneficial in avoiding the opportunities for bugs. For all of the voting systems that we consider here, the winner can be calculated in polynomial time. However, voting systems that use an IRV-type, round based structure would seem to be easier to implement than those that require graph algorithms such as Ranked Pairs and Schulze.
- Resistant to having too many “Schelling-dishonest scenarios”/Monotonic - We have seen in our work in [30] that, for any reasonable voting system and incentive system, there will always

---

Section 4.7.3. Thus, these incentive systems are not capable of distinguishing payouts between votes that rearrange lower ranked “irrelevant alternatives”.

<sup>21</sup>We briefly attempt to illustrate the differences in the attack space between WoodSIRV and Ranked Pairs, making use of examples. As both of these systems are Condorcet, if there is an “honest alternative”  $w$ , in both cases to change the result one must first manipulate enough votes to reverse at least one duel involving  $w$ . If reversing this duel results in a new Condorcet winner being selected, this same manipulation has the same result in the two systems. Or course, there are other possibilities for the profile of votes where one or the other voting system requires more manipulated votes to change the result. As an example of a situation where more manipulation is required to change the result in Ranked Pairs, consider the voting profile given by 2 votes for  $\{a > b > w\}$ , 5 votes for  $\{a > w > b\}$ , 8 votes for  $\{b > w > a\}$ , and 6 votes for  $\{w > a > b\}$ . Then  $w$  is the Condorcet winner, but already a single manipulation changing a  $\{a > w > b\}$  vote to  $\{a > b > w\}$  results in a Condorcet paradox where  $w$  actually receives the fewest first place votes and is eliminated under WoodSIRV. In order to change the result under Ranked Pairs, the attacker could take a few different strategies: 1) she could bribe/try to convince more votes to switch, for example, from  $\{a > w > b\}$  to  $\{a > b > w\}$  to strengthen the pair  $b > w$ , 2) she could bribe/try to convince more votes to switch, for example, from  $\{a > w > b\}$  to  $\{w > a > b\}$  to weaken the pair  $w > a$ , or 3) she could bribe/try to convince more votes to switch, for example, from  $\{w > a > b\}$  to  $\{w > b > a\}$  to weaken the pair  $a > b$ , or 4) some mixture of these strategies. While the changes in the first two strategies, to the degree that they do not successfully change the result from  $w$ , result in their voters incurring losses, the voters who follow the third strategy and vote  $\{w > b > a\}$  still get the maximum payoff if  $w$  remains the winner. As such, it might be more viable for an attacker to bribe or convince voters who would normally vote  $\{w > a > b\}$  to “micro-cheat” by voting  $\{w > b > a\}$  than it would to convince voters to manipulate their votes in other situations. On the other hand, consider the voting profile given by 5 votes for  $\{a > b > w\}$ , 2 votes for  $\{a > w > b\}$ , 5 votes for  $\{b > w > a\}$ , 7 votes for  $\{w > a > b\}$ , and 2 votes for  $\{w > b > a\}$ . Then an attacker could change the outcome under Ranked Pairs from voting for  $w$  as a Condorcet winner to selecting  $b$  in a Condorcet paradox by convincing a  $\{w > a > b\}$  vote to vote for  $\{b > w > a\}$ . In this case, WoodSIRV would still select  $w$  as the winner. Indeed, an attacker that wanted to change the outcome under WoodSIRV would have several available strategies: 1) bribe/convince enough votes to switch their rankings of  $a$  versus  $b$  such that  $b$  becomes a Condorcet winner, 2) Convince/bribe some of the votes who place  $a$  first to place  $w$  or  $b$  first so that  $a$  is eliminated in the IRV step, and  $b$  wins the duel with  $w$ . For example, the attacker might bribe an additional voter to change a vote from  $\{a > b > w\}$  to  $\{b > a > w\}$ . Both of these strategies might involve transposing pairs that are not weighted by the incentive system if  $w$  wins. If the attacker can successfully follow the first strategy, she could change the result of the vote under any Condorcet system. However, following the second strategy either requires transposing pairs that involve  $w$ , and hence taking penalties if  $w$  ultimately wins, or transposing pairs of alternatives that are already ranked both ahead of  $w$ . Hence, this limits the ability of the attacker to corrupt voters to make such transpositions to the limited set of voters that already bury (the “honest choice”)  $w$  and hence can already be thought of as being part of an attack coalition anyway.

be some rare situations where a voter is incentivized to deviate from what we call “Schelling-honesty”, namely she is incentivized to cast a vote that does not actually reflect an order in which she thinks alternatives are likely to win. This work develops upon classic impossibility theorems in social choice theory such as Arrow’s Impossibility Theorem [5] and the Gibbard-Satterthwaite Theorem [32] [54]. However, in that work we saw that the situations in which different voting systems failed Schelling-honesty were more contrived and seemingly less likely to occur in practice for some voting and incentive system pairs than for others. Notably, when a voting system is monotonic, the types of failures of Schelling-honesty are more limited.

The following table summarizes how selected voting systems perform according to these criteria:

	<b>Clone Independent</b>	<b>Complexity/ Gas</b>	<b>Condorcet</b>	<b>Monotonic</b>	<b>Attack Resistance</b>
<b>Plurality</b>	No	Low	No	Yes	Bad
<b>Borda</b>	No	Low	No	Yes	Not great
<b>Instant-Runoff</b>	Yes as voting system Bias in incentive system	Medium	No	No	Better?
<b>WoodSIRV</b>	Yes as voting system Bias in incentive system	Medium+	Yes	No	Better+?
<b>Ranked Pairs</b>	Yes as voting system No known bias in incentives	High	Yes	Yes	Better?
<b>Schulze</b>	Yes as voting system No known bias in incentives	High	Yes	Yes	Better?

When jurors are presented with a binary choice, as is the case of most current applications of Kleros<sup>22</sup>, these voting systems are equivalent. Indeed, such voting between two options satisfies all of the “good” properties in this table. Hence these comparisons are only relevant when there are at least three possible outcomes.

Finally, note that while the design choices of Kleros are motivated by the particular challenges of being attack resistant in a setting without trusted authorities, the properties considered above are potentially relevant beyond blockchain applications to other crowdsourced platforms that, like Amazon’s Mechanical Turk [1], require an aggregation of user feedback with users who may provide incorrect or spam answers to minimize the effort required of them. Hence some of these ideas could potentially be used to improve the design of such systems in the spirit of [46].

### 4.7.3 Incentive System(s)

Users are incentivized to become Kleros jurors as this gives them the opportunity to receive a portion of the arbitration fees paid for the dispute, as discussed in Section 4.6. As part of the incentive system that encourages jurors to provide honest rulings, in addition to the potential to gain arbitration fees, jurors can also lose some of their PNK stake for rulings that are out-of-line with those of the

<sup>22</sup>In fact, all Kleros disputes also have the possible outcome that jurors vote “refuse to arbitrate”, see Section 4.1. However, it is expected that this option will rarely be voted for, so cases with two other possible outcomes are typically de facto binary with only those two outcomes having plausible chances of being voted for.

other jurors. These lost PNK stakes are then redistributed to other, more coherent, jurors as will be described below. Thus jurors are participating in a Schelling game similar to those described in Section 2.

In order to be drawn in a given court, users are required to stake a minimum amount of tokens, denoted by `min_stake`. Then, regardless of how a juror votes on a case, the number of tokens that she can lose from her stake per vote is limited to a fixed percentage of this minimum stake. This percentage will be denoted by  $\alpha$ . However, as was observed in Section 4.5.2 in the discussion of the “weight”, a single juror can be drawn multiple times for a given case, giving her more votes on this case. Then the maximum amount that the juror can lose as a result of her vote increases corresponding to this weight. Namely, the maximum amount of tokens that can be lost on a given case per juror is:

$$D = \alpha \cdot \text{min\_stake} \cdot \text{weight}.$$

Both the  $\alpha$  and the `min_stake` parameters are defined by the governance mechanism and can vary from one court to another.

#### 4.7.3.1 Current Token Redistribution Model

Currently, in parallel with the Plurality voting system that we described in Section 4.7.2.1 (in which particularly jurors do not provide ranking of the options other than a single vote), any juror that does not select the outcome  $w$  that wins the last appeal round loses her deposit  $D$ . Then jurors that do vote for  $w$  receive a payment of:

$$\frac{\text{ETH fees and lost deposits}}{\# \text{ jurors that vote for } w}.$$

These calculations (i.e. how many deposits were lost and how many jurors voted for  $w$ ) are done on a round-by-round basis.

#### 4.7.3.2 Future Token Redistribution Model

Take  $w$  to be the option that wins via the vote aggregation methods described in Section 4.7.2. We speak about jurors voting “coherently” if they agree with the ultimate vote outcome; while being coherent is an all or nothing property for binary decisions, for non-binary decisions jurors’ votes can be more or less “coherent”. The goal is to incentivize users to place outcomes they believe to be “honest” high in their lists following the motivations of Section 4.7.2. Conversely, one also wants to strongly penalize a juror who has placed the winning choice  $w$  far down on her list. One option would be to have jurors lose:

$$\frac{\# \text{ options ranked above (or equal to) } w}{\# \text{ total options} - 1} D \tag{1}$$

to be redistributed between other jurors based on their coherence. However, in this framework, an attacker that ranks a malicious choice first and  $w$  second will risk relatively little of her deposit.

In equation 1, one rewards the jurors for the number of options  $a_i$  that they correctly place below the winner  $w$  with each  $a_i$  given the same weight. Alternatively, one can give extra weight for rewards and penalties for options  $a_i$  for which the margin of pairwise votes between  $w$  and  $a_i$  is particularly close. This is in the spirit that *narrowly* failed attacks should be particularly expensive, which is a common goal in the design of blockchain-based platforms [17]. If an attacker is attempting to commit

a bribing attack so that a would-be Condorcet winner  $w^*$  no longer wins, this requires a sufficient number of bribes so that at least some  $a_i$  defeats  $w^*$ . Hence at least one pair must pass from the honest winner winning to not, so in narrowly failed attack this pair will be weighted heavily.

Namely, one might take weights  $w(i) = w(a_i) \in [0, 1]$  for all  $a_i \neq w$ , such that  $\sum_{a_i \neq w} w(i) = 1$ . Then a voter  $USR$  loses:

$$D \sum_{a_j \neq w} \mathbf{1}_{USR \text{ voted } a_j \geq w} \cdot w(j)$$

from their deposit  $D$  and receives redistributions of the form:

$$\frac{\text{ETH fees and lost deposits}}{\sum_{USR_k \in \mathcal{V}} \sum_{a_j \neq w} \mathbf{1}_{USR_k \text{ voted } a_j < w} \cdot w(j)} \sum_{a_j \neq w} \mathbf{1}_{USR \text{ voted } a_j < w} \cdot w(j), \quad (2)$$

where  $\mathcal{V}$  is the set of voters in the same voting round as  $USR$ . In the following, we condense our notation by writing  $USR_k : a < b$  to indicate that the voter  $USR_k$  ranked the option  $a$  below the option  $b$ .

Note that, in a round with  $M$  voters, these rewards and penalties yield a net payoff for  $USR$  of

$$\begin{aligned} & \frac{\text{ETH fees} + \sum_k D \sum_{a_j \neq w} \mathbf{1}_{USR_k : a_j \geq w} w(j)}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{USR_k : a_j < w} w(j)} \sum_{a_j \neq w} \mathbf{1}_{USR : a_j < w} w(j) - D \sum_{a_j \neq w} \mathbf{1}_{USR : a_j \geq w} w(j) \\ &= \frac{\text{ETH fees} + MD}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{USR_k : a_j < w} w(j)} \sum_{a_j \neq w} \mathbf{1}_{USR : a_j < w} w(j) - D, \end{aligned}$$

with this last expression avoiding some redundant computation.

Under this payoff mechanism, regardless of which weight function  $w(i)$  is chosen, the lower jurors place the winning outcome, the larger the portion of the deposit they lose and the less arbitration fees they receive. Indeed, if a juror places the ultimate winning outcome last, she receives no arbitration fees and loses her entire deposit, which is split between other jurors in accordance with how high they ranked  $w^{23}$ . The above formulas do not involve a division by zero unless all voters in a given round place the winning outcome behind all other (non-zero weighted) alternatives, in which case all voters lose their deposit and do not receive a payout<sup>24</sup>.

For the choice of weight function, we have a series of tradeoffs and criteria similar to that of Section 4.7.2.2 for voting systems. Indeed, while we have divided this discussion into two sections for the sake of exposition, one for the voting system and one for the incentive system, the two choices can in some cases have interactions and should be considered together.

To illustrate the tradeoffs around a choice of weight function, we present a few of the weight functions that we have considered. In all cases, we will have an adjustable parameter  $\beta \geq 0$ , chosen via the governance mechanism that will control how concentrated the weights are on the closest pairs<sup>25,26</sup>

<sup>23</sup>The redistribution mechanism is inspired by the SchellingCoin, see Section 2, where jurors gain or lose tokens depending on whether their vote was consistent with the others jurors. Note that token redistribution mechanisms are still being actively researched and may further evolve in future versions.

<sup>24</sup>If at one level no one voted coherently, what to do with the amounts from that level can be determined by the governance procedure. See the descriptions on future work in Section 5 for further discussion on this point.

<sup>25</sup>Note that one might wish to choose different  $\beta$  in a way that depends on the dispute round as in early rounds, the number of jurors is small enough that which  $a_j$  are narrowly decided and hence which are weighted heavily is more variable, and hence to some degree arbitrary. For example, it can be helpful to take margins in the weight functions below drawn from the results of the last appeal round, when they will be the most representative of community sentiment.

<sup>26</sup>Further note that none of the weight functions considered below depend explicitly on which voting system we chose. In future work, one might incorporate into a weight function further information from the vote aggregation process, such as which round an alternative is eliminated in for an IRV-type system.

- Weight function 1 (constant weights):

$$w(i) = \frac{1}{\#\{a_i \in A : a_i \neq w\}},$$

- Weight function 2:

$$w(i) = \frac{\left(\frac{1}{|\text{margin of } a_i \text{ against } w|+1}\right)^\beta}{\sum_{a_j \neq w} \left(\frac{1}{|\text{margin of } a_j \text{ against } w|+1}\right)^\beta},$$

- Weight function 3:

$$w(i) = \frac{1 - \left(\frac{|\text{margin of } a_i \text{ against } w|}{\text{total number of votes}}\right)^\beta}{\sum_{a_j \neq w} 1 - \left(\frac{|\text{margin of } a_j \text{ against } w|}{\text{total number of votes}}\right)^\beta}$$

Then we evaluate these weight functions considering the following properties:

	<b>Winner First Max Payout</b>	<b>Unanimous No Weight</b>	<b>Concentration on Close Pairs</b>
<b>Weight Function 1</b>	Yes	No	No
<b>Weight Function 2</b>	Yes	No	Intermediate
<b>Weight Function 3</b>	Yes	Yes	Better

- Winner first gives maximum payout - Note that all of our weight functions have this property, indeed it follows from the redistribution structure given in equation 2. However, one could imagine payoff structures that depend not only on how voters rank the winning choice  $w$  compared to other choices, but also on the relative positions that voters give to two non-winning alternatives, i.e. whether a voter ranks  $a > b$  or  $b > a$  for some  $a, b \neq w$ . From a user experience perspective, we consider it important that users not be obliged to rank what they consider to be irrelevant choices, hence in all of our candidate systems here the relative rankings of such choices do not have weight. However, note that this choice fundamentally limits the ways in which a voting system+incentive system can be attack resistant. Indeed, by Arrow’s Impossibility Theorem, for any non-trivial voting system there will be situations where the relative rankings of “irrelevant alternatives” will affect the result [5]. For any incentive system that satisfies the Winner First, Maximum Payout property, a voter’s ranking of these “irrelevant alternatives” does not affect her payout, so the threshold to “micro-corruption” attacks, for example bribing such a juror to invert these seemingly less relevant alternatives, may be lower. See our comments in Section 4.7.2.2 on why we think this effect is mitigated, while nonetheless present, when the voting system used satisfies Later No Harm and Later No Help.

- Unanimous pairs given no weight - Generally, the weight functions above (with the exception of the constant weight function) are such that  $a_i$  is given more weight than  $a_j$  if the duel between  $w$  and  $a_i$  is closer than that between  $w$  and  $a_j$ . Taken to the extreme, one might hope that duels that are unanimous are given no weight, to not dilute the effect of the payoff system on incentivizing voters to make honest rankings on pairs that are more likely to matter. Indeed, for weight functions that have this property, if a non-binary case is “defacto binary”, i.e. there are two alternatives  $a$  and  $b$  such that all voters rank  $a$  and  $b$  first and second in some order, then (assuming that the voting system actually selects  $a$  or  $b$  as the winner, which would be the case for any of the voting systems considered in 4.7.2.2) only the duel between  $a$  and  $b$  would receive weight in the incentive system and this situation reduces to the incentives one has analyzed in binary cases.
- Concentration of weight on closes pairs -This is related to, though somewhat more general than Unanimous, No Weight. As previous discussed, one wants to weight closer pairs more to increase attack resistance. However, one still wants to give some weight to other pairs to encourage voters to take their entire vote (or at least the rankings between all alternatives that they think have a meaningful chance of winning) seriously. For now at least our evaluations of different weights on this criterion are largely heuristic: for different values of  $\beta$  and common vote splits (e.g. 2-1, 5-2, or 12-3) how much weight is given to each alternative, and how the shape of these weight functions varying in  $\beta$  influences the flexibility of the governance process 4.12 to adjust.

In the following proposition, one sees that this payoff system can have good properties with respect to incentivizing jurors to rank candidates who are likely to win higher, corresponding to the objectives laid out in Section 4.7.2.

**Proposition 1.** *Consider the incentive system above with the constant weights, i.e. weighting function 1. Suppose that a given voter has a probabilistic prior for the outcome of the dispute resolution process, i.e. she estimates probabilities for the votes of other jurors and for the probabilities of each outcome to win possibly after appeals, such that:*

- *she believes that the votes of other jurors in her voting round are independent of her vote,*
- *she believes that the outcome is independent of her vote and the votes of other jurors in her voting round,*
- *she assigns to the possible outcomes  $a_1, \dots, a_n$  probabilities  $prob(a_1), \dots, prob(a_n)$  of ultimately winning.*

*Then a weakly dominant strategy for this juror is to provide a (strict) ranking of the outcomes  $a_j$  from highest to lowest by their chance of winning  $prob(a_j)$ .*

See Appendix A for a proof of this result. Note that the perspective of a juror believing that her vote will not change the ultimate outcome can be justified in our setting if jurors believe that incorrect outcomes are likely to be appealed.

After Kleros has reached a decision, tokens are unfrozen and redistributed among jurors. An example of token redistribution is shown in Figure 10. Note that jurors could fail to reveal their vote. To disincentivize this behaviour, the penalty for not revealing one’s vote is at least as large as the penalty for voting incoherently. This incentivizes jurors to always reveal their vote. In case of appeal, arbitration fees and tokens are redistributed at each level according to the result of the final appeal.

When there is no attack, parties are incentivized to vote what they think, other parties think, other parties think... is honest and fair. In Kleros, the Schelling Point is honesty and fairness. One could

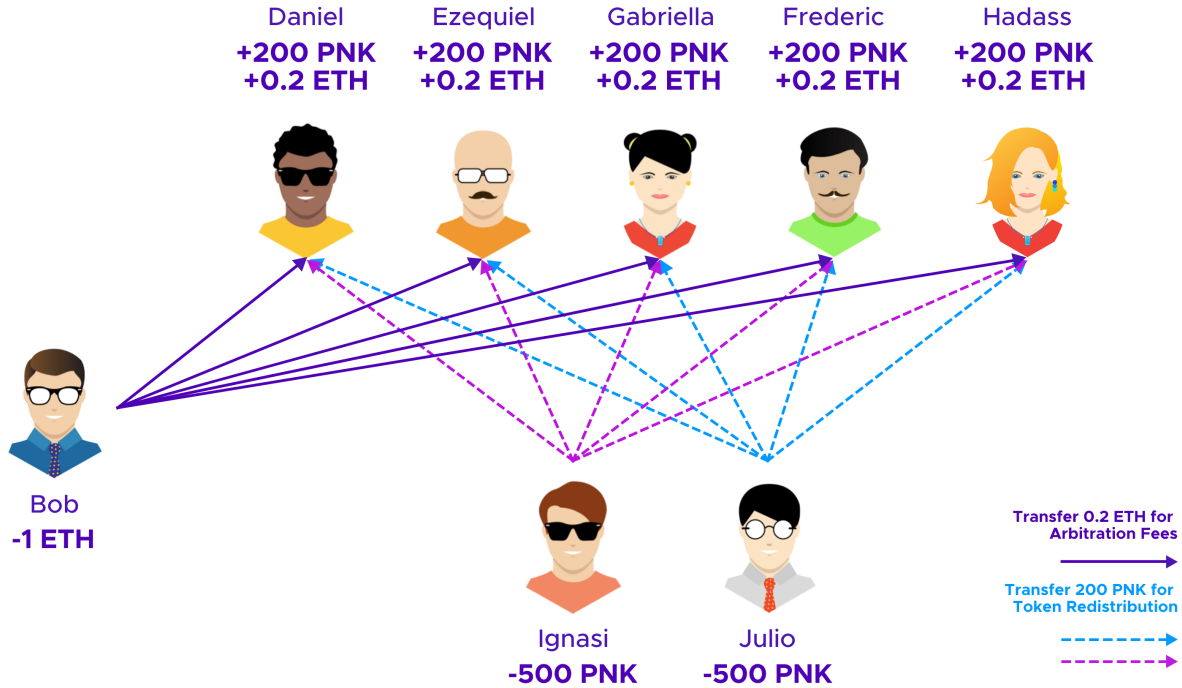


Figure 10: Seven jurors have a binary choice between ruling on behalf of Alice or on behalf of Bob. Tokens are redistributed from jurors who voted incoherently to jurors who voted coherently. Bob lost the dispute and pays the arbitration fees. The other deposits are refunded.

argue that those decisions being subjective (for example, compared to a SchellingCoin mechanism for a prediction market), no Schelling Point would arise. In [55], the informal experiments run by Thomas Schelling showed that in most situations a Schelling Point plebiscited by all parties does not exist. But Schelling found that some options were more likely to be chosen than others. Therefore even if a particularly obvious option does not exist, some options will be perceived as more likely to be chosen by others parties and will effectively be chosen. We cannot expect jurors to be right 100% of the time. No dispute resolution procedure could ever achieve that. Some times, honest jurors will lose coins. But as long as overall they lose less value than what they win as arbitration fees and as coins for other incoherent parties, the system will work<sup>27</sup>.

**Remark 1.** Above we saw that the redistribution of arbitration fees and lost deposits is handled by round. Note that if a given voter then knows or suspects that other voters in her round have voted “incorrectly”, this gives her more of an incentive to vote honestly. In the extreme, a single juror that agrees with the final outcome in a round where every other juror disagreed would receive all the arbitration fees and lost deposits for that round. We call this phenomenon the “lone voice of reason” effect. We will note further implications of this effect below.

<sup>27</sup>Indeed, note that, so far, this idea has largely worked as expected in experiments such as those considered in [29] as well as in practical applications such as those of [34] and [49].



#### 4.7.4 Conclusions

Weighing these considerations, for typical cases involving discrete multiple choices, a good choice of voting and incentive system seems to be:

- the voting system WoodSIRV (i.e., one uses the rankings provided to simulate a series of voting rounds, eliminating the alternative with the fewest first place votes in each round, as in IRV, but before each round one checks if there is a Condorcet winner among the remaining alternatives)<sup>28</sup>, and
- weight function 3, namely

$$w(i) = \frac{1 - \left( \frac{|\text{margin of } a_i \text{ against } w|}{\text{total number of votes}} \right)^\beta}{\sum_{a_j \neq w} 1 - \left( \frac{|\text{margin of } a_j \text{ against } w|}{\text{total number of votes}} \right)^\beta},$$

where  $\beta$  is an adjustable parameter (that can take an arbitrary positive value). Particularly, the weight function should be calculated based on the margins from the last voting round (in order to be based on the largest, most statistically significant sample) when calculating the rewards for voters in earlier rounds.

Nonetheless, to the degree that certain applications might want to prioritize different properties in the tradeoffs we have considered above, the modular structure of future versions of the court can allow those applications to use different voting and incentive systems, to the degree that those alternative models have been approved by the governance process, see Section 4.12.

#### 4.7.5 Parameterization of Arbitration Fees

Suppose we have a round of  $M$  jurors. One must choose  $f$ , the per juror average juror fee (i.e. so that the entire round requires  $M \cdot f$  in arbitration fees), as well as the parameters  $\alpha$  and  $\text{min\_stake}$  described above that give the maximum deposit per juror vote that can be lost:  $D = \alpha \cdot \text{min\_stake}$ .

Again, denote the ultimate winning outcome by  $w$ . In the model of Section 4.7.3.2, imagine an honest juror that takes time and effort valued at  $e$  and as a result has:

$$E \left[ \frac{\sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}:a_j < w} w(j)}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k:a_j < w} w(j)} \mid \exists k : \mathcal{USR}_k \text{ does not put } w \text{ last} \right] \geq \frac{1}{M},$$

namely that, on average, her return is at least as high as the average juror return.

Again, we can express  $\mathcal{USR}_j$ 's net return as

$$\frac{Mf + \sum_k D \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k:a_j \geq w} w(j)}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k:a_j < w} w(j)} \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}:a_j < w} w(j) - D \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}:a_j \geq w} w(j)$$

---

<sup>28</sup>An important technical point in the choice of a voting system is how to handle ties. For WoodSIRV, a particularly simple choice of tiebreaker from the perspective of minimizing code complexity is to require that an alternative *strictly* win all of its duels to be considered a Condorcet winner, and when eliminating options in last place to eliminate all alternatives that are tied for last place. Finally, one can return “refuse to arbitrate”, similarly to the current mechanism for handling tied options, if this process results in all options being eliminated. However, as such an implementation deviates slightly from WoodSIRV as it is typically defined, it may fail some properties in edge cases that normal WoodSIRV satisfies (e.g. failures of Clone Independence as a voting system could occur if all members of a set of clones were tied for last in a given voting round).

$$= \frac{Mf + MD}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j < w} w(j)} \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}: a_j < w} w(j) - D.$$

Denote

$$\pi_* = \text{prob} \left( \begin{array}{l} \mathcal{USR}_k \text{ puts} \\ w \text{ last } \forall k \end{array} \right)$$

Then, we can calculate the expected value of this honest strategy, based on the payoffs above:

$$\begin{aligned} E[\text{honest}] &= (1 - \pi_*) E \left[ \text{honest} \mid \begin{array}{l} \exists k : \mathcal{USR}_k \text{ does} \\ \text{not put } w \text{ last} \end{array} \right] - \pi_*(D + e) \\ &= (1 - \pi_*) E \left[ \frac{Mf + MD}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j < w} w(j)} \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}: a_j < w} w(j) - D - e \mid \begin{array}{l} \exists k : \mathcal{USR}_k \text{ does} \\ \text{not put } w \text{ last} \end{array} \right] - \pi_*(D + e) \\ &\geq (1 - \pi_*)(f - e) - \pi_*(D + e). \end{aligned}$$

On the other hand, noting that the special case where all voters in a round rank  $w$  last corresponds to the minimum payoff, a “lazy” strategy adopted by the voter  $\mathcal{USR}_l$  that does not expend any effort in the case has expected value:

$$\begin{aligned} E[\text{lazy}] &= (1 - \pi_*) E \left[ \frac{Mf + MD}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j < w} w(j)} \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}: a_j < w} w(j) - D \mid \begin{array}{l} \exists k : \mathcal{USR}_k \text{ does} \\ \text{not put } w \text{ last} \end{array} \right] - \pi_* D \\ &= (1 - \pi_*) M(f + D) E \left[ \frac{\sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_l: a_j < w} w(j)}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j < w} w(j)} \mid \begin{array}{l} \exists k : \mathcal{USR}_k \text{ does} \\ \text{not put } w \text{ last} \end{array} \right] - D. \end{aligned}$$

As long as the best available “lazy” strategy has a lower average payoff than that received by the average juror, i.e. as long as:

$$E \left[ \frac{\sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_l: a_j < w} w(j)}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j < w} w(j)} \mid \begin{array}{l} \exists k : \mathcal{USR}_k \text{ does} \\ \text{not put } w \text{ last} \end{array} \right] \leq \frac{1}{M},$$

it is possible to choose  $D$  sufficiently large such that the lazy strategy has a negative expected return. Then, one should choose  $f$  and  $D$  such that:

$$(1 - \pi_*)(f - e) - \pi_*(D + e) > 0 > (1 - \pi_*) M(f + D) E \left[ \frac{\sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_l: a_j < w} w(j)}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j < w} w(j)} \mid \begin{array}{l} \exists k : \mathcal{USR}_k \text{ does} \\ \text{not put } w \text{ last} \end{array} \right] - D.$$

Here one can view  $\pi_* = \text{prob} \left( \begin{array}{l} \mathcal{USR}_k \text{ puts} \\ w \text{ last } \forall k \end{array} \right)$  and  $E \left[ \frac{\sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_l: a_j < w} w(j)}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j < w} w(j)} \mid \begin{array}{l} \exists k : \mathcal{USR}_k \text{ does} \\ \text{not put } w \text{ last} \end{array} \right]$  as quantities that can be observed empirically in a given court over time.

We can perform a similar analysis under the model of Section 4.7.3.1. We suppose that the honest strategy determines the winning answer with probability  $p$ . Suppose that all jurors other than  $\mathcal{USR}$  adopt the honest strategy, so if we consider the jurors as choosing their answers independently, the number of other jurors who vote for the ultimate winning answer is binomial  $X \sim \text{Binom}(M - 1, p)$ . Suppose a “lazy” juror that does not evaluate the case can choose the right answer with probability  $t \in [0, p]$  (for example because the court tends to side with a contractor over a business owner in a  $t$  proportion of cases). Then:

$$E[\text{honest}] = pE \left[ \frac{Mf + (M - X - 1)D}{X + 1} \right] + (1 - p)(-D) - e = (f + D)(1 - (1 - p)^M) - D - e$$

and

$$E[\text{lazy}] = tE\left[\frac{Mf + (M - X - 1)D}{X + 1}\right] + (1 - t)(-D) = (f + D)\frac{t}{p}(1 - (1 - p)^M) - D.$$

Namely, in this case, the above constraints become:

$$(f + D)(1 - (1 - p)^M) - D - e > 0 > (f + D)\frac{t}{p}(1 - (1 - p)^M) - D.$$

If these are satisfied, the honest strategy is Bayesian-Nash incentive compatible. Note that in the special case of the model of Section 4.7.3.2 where jurors make a binary choice hence  $w(j) = 1$ , the two models are the same. We calculate the probability that no juror votes for the ultimate winner as  $(1 - p)^M$ , and

$$E\left[\frac{\sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_i: a_j < w} w(j)}{\sum_k \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j < w} w(j)} \mid \exists k : \mathcal{USR}_k \text{ does not put } w \text{ last}\right] = \frac{t}{1 - (1 - p)^M} E\left[\frac{1}{X + 1}\right] = \frac{t}{Mp}.$$

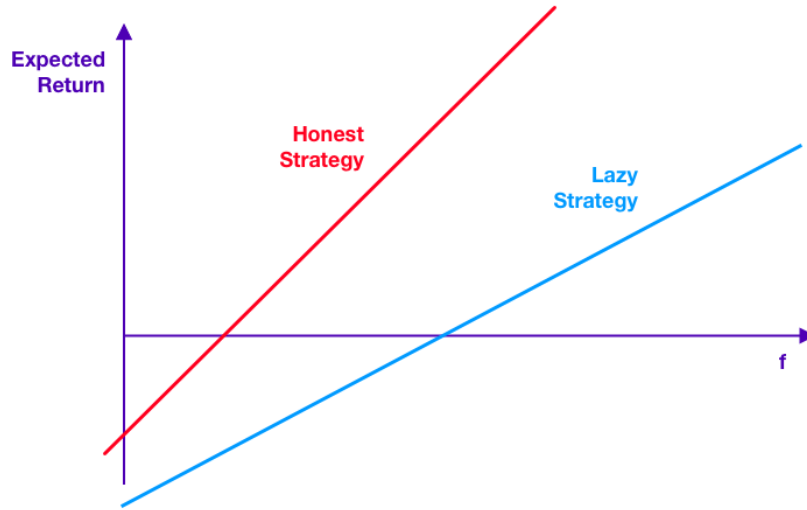


Figure 11: The parameters  $f$  and  $D$  should be chosen such that the strategy of performing effort to review the case and voting honestly has a positive expected value and such that the lazy strategies of voting randomly or always making the same vote have a negative expected return. Here the choice of  $D$  affects the positions of these curves and should be chosen such that there is an acceptable range from which to choose  $f$ .

Ultimately, the values  $f$  and  $D$  are chosen by the governance mechanism, see Section 4.12. Particularly, in order to follow this reasoning, the community must estimate values such as  $e$ <sup>29</sup>. Hence, if adequate data is available to justify using a more sophisticated model for how jurors' efforts vary over the population, the governance votes are capable of following more subtle versions of this argument.

<sup>29</sup>Note that one would want the different cases being considered by a given court to generally require similar levels of effort so that jurors do not apply the "honest" strategy on easier cases and the "lazy" strategy on difficult cases. Nevertheless, even if the cases in the court are of roughly equal difficulty, some jurors will need to exert less effort than others to come to honest rulings.

### 4.7.6 Enforcing Standards of Juror Behaviour

Several aspects of the behaviour that is expected from Kleros jurors, as laid out in court policies, are subjective, and smart contracts are not capable of directly enforcing such behaviour. For example, while court policies may require jurors to provide sensible justifications with their votes, it can be difficult for an automated process to distinguish between a meaningful text provided as a justification and a random string. Similarly, it is difficult to create fully automated processes to disincentivize jurors from prematurely revealing their vote during the vote period, see Section 4.9.1.

Hence, future versions of Kleros will make use of the ability of Kleros itself to render judgments on subjective questions, by having a special court in which rulings are made on whether given jurors violated some court policy. Namely, during a given period that includes and extends beyond the voting period, a juror Alice to a given dispute can be challenged as having violated some juror requirement with respect to her vote in the case. This challenge is decided via a separate dispute in a “Process Court”.

Regardless, of whether the challenge is successful or not, the challenger pays any arbitration fees that might be required to raise a dispute in the Process court. The challenger is incentivized to take on these costs by the possibility to win a deposit from Alice’s minstake.

Namely, for each court, there is a parameter  $\gamma$  chosen by governance 4.12. If jurors rule that Alice violated court policies, the challenger receives an amount of:

$$\gamma \cdot \text{min\_stake} \cdot \text{weight}$$

from Alice’s PNK stake. If the challenge is unsuccessful, Alice receives a deposit that is lost by the challenger.

## 4.8 Appeals

If, after the jury has reached a decision, a party is not satisfied (because she thinks the result was unfair), she can appeal and have the dispute ruled again. Each new appeal instance will have twice the previous number of jurors plus one. Hence, these larger panels will generally have a progressively higher likelihood of being statistically representative of the broader juror community. Due to the increased number of jurors, appeal fees must be paid ( $\text{appeal\_fees} = \text{new\_amount\_jurors} \cdot \text{average\_fee\_per\_juror}$ ).

The number of jurors increases exponentially as one appeals; hence arbitration fees also rise exponentially with the number of appeals. This means that, in most cases, parties won’t appeal, or will only appeal a moderate amount of times. Hence, via the appeal mechanism, Kleros manages to avoid the unnecessary duplication of effort and high costs that would be required by having a very large number of jurors consider every case while nonetheless the possibility of appealing a high number of times provides a defense against an attacker bribing the jurors. See Section 4.11 for a further discussion of this point.

### 4.8.1 Appeal Courts

In case of appeal, a case can be considered by a larger panel of jurors in the same court in which it had previously been considered or, in some cases, it can be raised to a the “parent court” of this court, see the tree of courts discussed in Section 4.3.

Note that each “child court” typically will require additional specialized skills compared to its “parent court”. Hence, when a case progresses to a different court, it goes to a court where jurors are somewhat more general and may be less specialized in the particular relevant use case, resulting in jurors in the parent court requiring more effort to come to an appropriate ruling and/or obtaining

this “honest” ruling with somewhat less consistency. On the other hand, more jurors will be staked in this parent court; hence it has a greater resistance to various attacks, see Section 4.11.

Hence, the different courts exhibit a tradeoff between specialization and attack resistance. In order for cases to typically benefit from the skills of specialized jurors in smaller courts, while still being attack resistant, court policies instruct jurors in appeal courts to look for evidence of attacks. For example, the General Court policy includes the instructions:

*“When considering an appeal of a case that has originated in a lower court, jurors should consider whether 1) evaluating the case requires specialized skills which jurors in the appellate court cannot be expected to have (ex: evaluating the quality of an English to Korean translation when knowledge of Korean is not a requirement of the appellate court) and 2) whether there is evidence that an attack was performed against this case in the lower court (ex: bribes, p+epsilon attacks, 51% attacks, etc). If there is no evidence of an attack AND appellate court jurors cannot be reasonably expected to have the required skills to independently evaluate the case, jurors should vote to uphold the lower court ruling”* [39].

This structure of appeal courts focusing on process questions is similar to the role of appeal courts in federal court systems such as that of the United States [4]. This process of appealing from a child court to its parent can continue all the way until a case is in the General Court.

In the current model, an appealed case “jumps” to a higher court, namely it is considered in the parent court of the court where it was previously considered if:

$$\begin{array}{l} \text{number of jurors} \\ \text{in previous round} \end{array} > \text{threshold},$$

where the threshold here is court dependent and is set via the governance process described in Section 4.12. As the number of jurors in each round is double the number of jurors in the previous round plus one, this structure essentially determines how many appeal rounds a case should be considered in each court before moving on to subsequent appeal courts<sup>30</sup>. If the number of jurors in a round exceeds the value of this threshold in the General Court, then any further appeal would result in a forking vote, see Section 4.10.

In future versions of the court, in addition to the number of jurors surpassing a cutoff being sufficient to send appeals to higher courts, there will also be a condition that triggers a jump to a higher court if:

$$\text{prob} \left( \begin{array}{l|l} \text{previous round} & \text{average juror would} \\ \text{result or} & \text{vote current winning choice} \\ \text{more extreme} & \text{first with } < \frac{1}{2} \text{ probability} \end{array} \right) \leq \text{threshold}.$$

This probability can be estimated via binomial tail bounds. Again, the threshold here is an adjustable parameter that is set via the governance process described in Section 4.12. This condition allows one to avoid appeals in a court where the previous vote of the jurors in that court is so overwhelming that it is statistically extremely likely that the current winning alternative already represents the general opinion of jurors in that court.

---

<sup>30</sup>Note during the governance process to update these thresholds, one can consider the number of jurors and diversity of the stakes in the different courts. Particularly, one can take lower values for the threshold in more specialized courts where there are fewer distinct jurors staked. Namely, one can choose the thresholds such that when further appeals in that court are unlikely to draw new jurors to provide a fresh analysis of a case, the case will proceed to the parent court where it will be considered by a broader pool of jurors.

## 4.8.2 Appeal Fee Models and Crowdfunding

In the current version of Kleros, arbitrable contracts specify their models for gathering the required appeal fees and determine the consequences when fees are not paid. These models present various tradeoffs that we will discuss below. A first, basic choice is the following:

- In the case of appeal, the appellant must pay any required appeal fees. This has the advantage that the party that won the previous round cannot be caused to lose without a further ruling by the jurors merely because she did not pay appeal fees. On the other hand, this has the disadvantage that parties who ultimately win their appeals do not have their appeal fees reimbursed, hence it may only be worth appealing relatively high value disputes.
- Alternatively, each side can be required to pay sufficient appeal fees to cover necessary costs in case they lose the appeal. Then, if only one option is (fully) funded, then that option would win by default without an additional appeal round. To mitigate the issue that a party that won a previous round can lose merely as a result of not paying adequate fees, we require additional stake to be paid by parties funding an appeal beyond the fees required to pay jurors for the subsequent round, where the option that won the previous round might require less stake. Then this stake is used to incentivize “fee funders”. Namely one encourages third parties to pay the appeal fees for options, particularly those of the previous round winner, in exchange for the possibility of winning the stake of the other side. This structure is similar to litigation funding [52]. Note that fees can be funded collectively by a “crowd”, as we will describe below<sup>31</sup>.

In the next version of Kleros, this choice over who can pay appeal fees and in what circumstances they can do so will be incorporated into the arbitrator contract. However, the modular structure of this version of the court will allow parties creating disputes to choose between different models for the structure of appeal fees, to the degree that those alternative models have been approved by the governance process, see Section 4.12.

In the rest of this section we will describe mechanisms, in which third party “funders” are incentivized to cover appeal fees on behalf of parties to a dispute. We will provide several models for doing so, again each with its own tradeoffs. In designing these models, one has the following constraints:

- at least the first funder in an appeal period has to stake on a given, single choice so that if the fees of no other choice are paid then there is no dispute and that choice becomes the default winner.
- at least one funder must be wrong per round so that they can cover the appeal fees for that round.

We use the following notation for this section:

- $x$  is the total fees required by the jurors in the subsequent appeal
- $s_{a_i}$  is the additional stake required beyond arbitration fees in order to fund the outcome  $a_i$ ; namely, the total deposit that must be made to fund  $a_i$  is  $x + s_{a_i}$ <sup>32</sup>.

---

<sup>31</sup>Other mechanisms designed to protect less well-financed parties that may not be able to pay large appeal fees are also possible and are an active subject of research. For example, parties to disputes can participate in a collective “appeal fee insurance”. Here the parties would deposit amounts greater than the required first round juror fee when initially creating their dispute. Then the loser of the case does not receive back this difference, rather it goes into a pool of money that is used to pay appeal fees of parties that had won the previous round when required. These different models, fee litigation funding and fee insurance, each have their own tradeoffs, though they can be used together.

<sup>32</sup>It is possible that during the course of an appeal, the governance mechanism for a court will change its required arbitration fees. A variety of options for how to handle this, from requiring sides that had previously been fully funded to contribute more to requiring the difference from the remaining party, are possible.

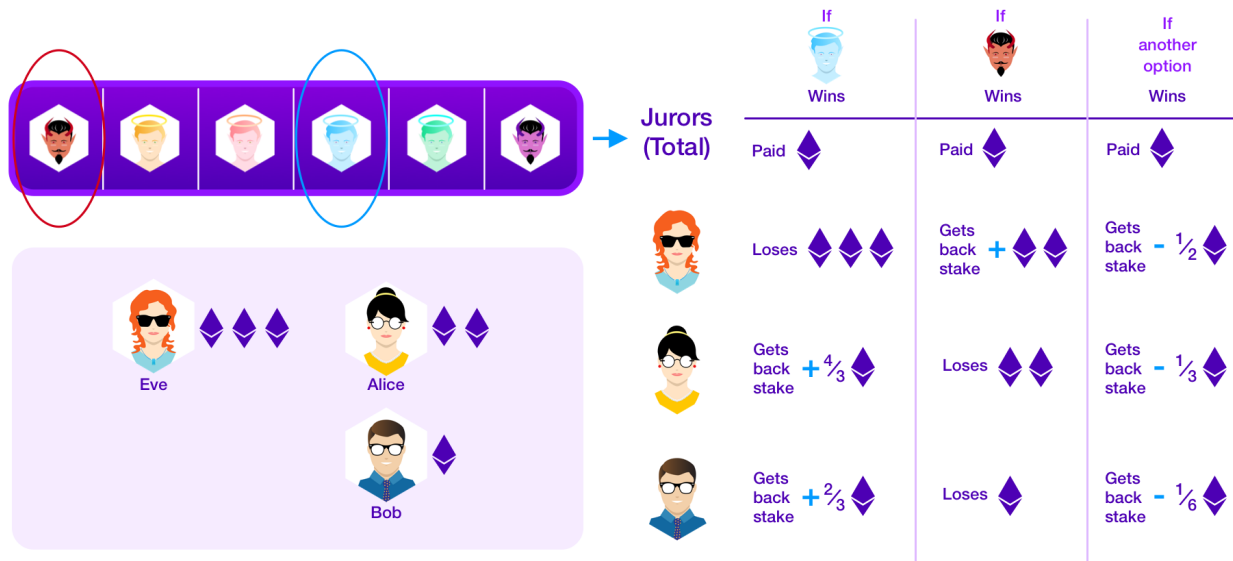


Figure 12: If Eve funds the appeal fees of a dishonest choice, Alice and Bob can collectively crowdfund the appeal fees of another (honest) choice. Crowdfunders are incentivized to participate because they can win a stake that is paid from the appeal fees of the opposing side beyond what is required to pay the jurors. Any situation where only two options are funded resembles this schema, regardless of which of the crowdfunding models below is used.

Different designs can be adopted:

- In order to fund an option  $a$ , one must deposit  $x + s_a$ . Then once two options are funded, a dispute is raised. If some other option ultimately wins, the two funders share the arbitration fees for their round, but otherwise their deposits are returned. We will see that this model has reasonably good resistance to the “clone funding” grief described in Section 4.11. However if there are multiple options that could be considered to be honest, this model does not give the possibility to the funder to hedge by funding them together.
- During the appeal period, users can pay  $x + s_{a_i}$  for the choice  $a_i$  to be funded. Any choice that is not funded is eliminated from juror consideration in all future rounds. This has the advantage that cases are likely to binarize themselves after a small number of rounds, removing complications coming from having multiple choices. On the other hand, winnowing possible outcomes primarily through requirements to pay fees creates the risk that no good option is funded in a round and jurors are presented with nonsensical disputes. Moreover, an approach to winnowing the set of options that depends so heavily on fees paid could be seen as plutocratic. Also this has the disadvantage that it is incompatible with the current standard for arbitrator and arbitrable contracts [43].
- Suppose that option  $b$  won the previous appeal round. Then an appeal is called if for some option  $a \neq b$  and some set of options  $S$  that does not contain  $a$  but may contain  $b$ ,  $x + s_a$  is staked on behalf of  $a$  by one or more funders and  $\gamma(\#S)x + \sum_{a_i \in S} s_{a_i}$  is staked on behalf of  $S$  by one or more funders, where  $\gamma(\#S)$  is an increasing function satisfying  $1 \leq \gamma(\#S) \leq \#S$  for all values  $\#S$ <sup>33</sup>. If  $a$  wins, the funders of  $a$  receive back their deposits plus  $(\gamma(\#S) - 1)x + \sum_{a_i \in S} s_{a_i}$  proportionally based on their contributions. If an option in  $S$  wins, then the funders who contributed to the fees

<sup>33</sup>Note that a compatible set  $S$  can be found if it exists by taking all options  $a_i$  for which at least  $s_{a_i}$  has been raised.

of  $S$  receive  $s_a$  proportionally based on their contributions. If an option outside of  $\{a\} \cup S$  wins, then the side insuring  $a$  receives back her deposit minus  $\frac{x}{\#S+1}$  and the funders of  $S$  receives back their deposits minus  $\frac{x\#S}{\#S+1}$ . If there are multiple options that could be considered to be honest, this model gives the possibility to the funder to hedge by funding them together. However, for the simplest choice of  $\gamma$ ,  $\gamma(\#S) = 1$ , this model has worse resistance to the “clone funding” grief described in Section 4.11.6 than the previous model. We will see in Proposition 7 that if one chooses  $\gamma(k) = \frac{k+1}{2}$ , this model has comparable resistance to clone funding as the first model, while maintaining its good properties with respect to allowing hedging for funders, at the expense of additional complexity.

- Finally, we consider a slight variation on the previous model, where instead of a single set of options  $S$  being chosen, different crowdfunders can stake on different sets of options. We take  $\gamma$  as above, and then an appeal is triggered if for some option  $a \neq b$ ,  $x + s_a$  is staked on behalf of  $a$  by one or more funders and other funders stake on set of options not containing  $a$  such that

$$\sum_{\substack{S \subseteq A \\ a \notin S}} \frac{\text{Amount funded for } S}{\gamma(\#S)x + \sum_{c \in S} s_c} = 1. \quad (3)$$

Then returns for various participants are given by

$$\sum_{\substack{S \subseteq A \\ a \notin S}} \frac{\text{Amount funded for } S}{\gamma(\#S)x + \sum_{c \in S} s_c} \cdot \begin{matrix} \text{Return under prev. model if} \\ S \text{ had been funded alone} \end{matrix}.$$

Note, in particular, the amount paid to jurors is  $x$  independently of which sets are funded. Moreover, when  $s_c = s$  is fixed for all  $c \neq b$ , the sums above can be efficiently computed as sums over the participating crowdfunders, for example equation 3 becomes

$$\sum_{\text{crowdfunder } \mathcal{U}SR_j} \frac{\text{Amount funded by } \mathcal{U}SR_j \text{ for } S_j}{\gamma(\#S_j)x + \#S_j \cdot s + (s - s_b) \cdot \mathbf{1}_{b \in S_j}} = 1.$$

We will see that this model has many of the advantages of the model where a single set  $S$  is funded; however, it provides more flexibility to crowdfunders to stake on different sets and nonetheless collectively fund an appeal.

In all of these models, if insufficient funds are raised for a given option or group of options to trigger an appeal, those funds should be returned to their contributors. This encourages funders to participate in this process without taking unnecessary risks on whether an option will manage to be fully funded.

**Remark 2.** *One notes that in our models above we allow the possibility for several choices to be funded together only on one “side” of the dispute. Ideally, one could place any elaborate “bet” on either side of the dispute, funding some combination of outcomes to hedge as necessary. However, as mentioned in our constraints above, the first “side” that is funded must indicate a single option that becomes the default option if no other fees are paid. One could imagine attempts to adapt, such as taking the highest ranked choice from the previous round among a set of options that are funded together to be the default choice. However, this is vulnerable to attacks where a malicious party that has managed to corrupt a previous round can ensure that all high ranking choices are malicious and then fund one high ranking choice together with enough honest choices that it will not be viable for honest parties to fund another side to provoke a dispute.*



**Remark 3.** *This crowdfunding mechanism does not completely resolve the issue of making appeal fees accessible. While (different) funders can cover the fees of different sides of a given dispute, they are ultimately playing a negative sum game against each other as the funders collectively must pay the arbitration fees for the appeal round. In particular, funding two sides of a dispute would typically only be possible in cases where funders have very different priors for the chances of the victory of different outcomes. Indeed, there will generally be ranges of the funders' priors in which the case is not sufficiently clear-cut to fund any party. This issue becomes more acute as the number of possible outcomes grows; with many options, an attack could push a dispute to some clearly false option without there being a single honest option with sufficient chances of victory to be funded in appeal. This is particularly true in the presence of clones that may drive down the probability that any given honest answer wins. Hence, when designing an arbitrable contract, one should consider combining crowdfunding with pre-ruling appeal fee insurance as well as putting in careful reflection on the possible outcomes proposed, avoiding unnecessary clones.*

**4.8.2.1 Funding Sets of Choices Together**

We briefly show a few results in the context of allowing multiple choices to be funded together, i.e. the third and fourth crowdfunding models discussed above. Particularly, we see that participants funding additional choices does not harm actors who had already (partially) funded. We first see that it can make sense for funders on the side of a set of choices to contribute to several different choices as this can maintain positive expected value of return while reducing variability.

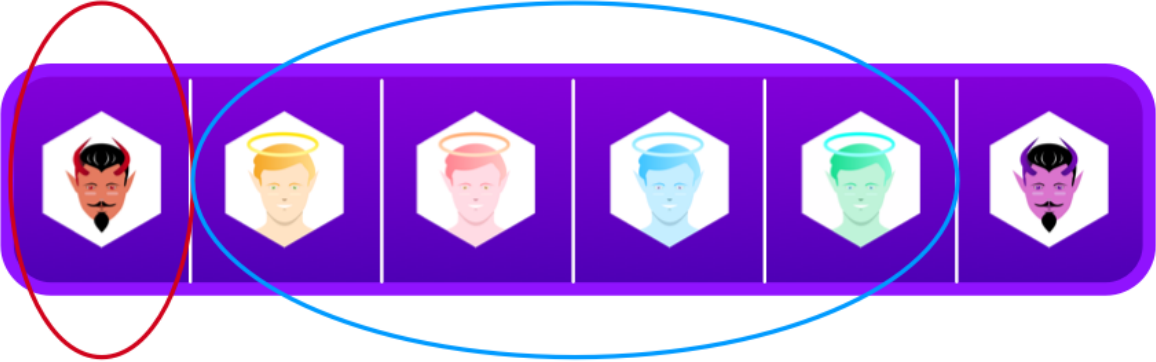


Figure 13: When one employs the third or fourth crowdfunding model above, funders can finance the fees of a set of options. Here, in contrast to the situation of Figure 12, crowdfunders stake on the honest options collectively and win if any of those options ultimately wins the dispute. This requires a larger contribution from the crowdfunders than if they had only funded a single option, however generally the additional cost of adding an option is sub-additive, with the degree of sub-additivity controlled by the choice of  $\gamma$ . This allows for a form of hedging that can preserve the effectiveness of crowdfunding even in the presence of clone options.

**Proposition 2.** *Suppose the function  $\gamma$  has the property that  $\gamma(n) + \gamma(1) \geq \gamma(n + 1)$ , for all  $n \in \mathbb{N}$ . Then if funding the set of outcomes  $\mathcal{S}$  has a non-negative expected value and funding the outcome  $a_j$  individually has a non-negative expected value, then funding  $\mathcal{S} \cup \{a_j\}$  collectively has a non-negative expected value.*

*Proof.* Suppose that  $a$  is the option that has been funded by the opposing side. Take  $\#S = n$ . Then

$$E[S \cup \{a_j\}] \geq E[S] - E[a_j] \Leftrightarrow$$

$$p_a x (-\gamma(n+1) + \gamma(n) + \gamma(1)) + (1-p_a) \left( \frac{nx}{n+1} + \frac{x}{2} - \frac{(n+1)x}{n+2} \right) + \sum_i p_{a_i} \left( \frac{(n+1)x}{n+2} - \frac{nx}{n+1} \right) + p_{a_j} \left( \frac{(n+1)x}{n+2} - \frac{x}{2} \right) \geq 0.$$

The first term is non-negative by our assumption on  $\gamma$ . The other terms are clearly non-negative as  $n \geq 0$ . □

**Proposition 3.** *Suppose that funding the option  $a$  has a non-negative expected value if the opposing side has funded  $a_j$  for any option  $a_j \in \mathcal{S}$ . Then funding  $a$  has a non-negative expected value if the opposing side has funded  $\mathcal{S}$  collectively. Moreover, in the fourth model above, if funding the option  $a$  has a non-negative expected value when the opposing side has funded  $S_j$  for any  $S_j$  such that  $a \notin S_j$ , then funding the option  $a$  has a non-negative expected value for any choices of stakes for which equation 3 holds.*

*Proof.* Let  $n = \#\mathcal{S}$ . Then for  $n \geq 1$ ,

$$E[a] = p_a \left( (\gamma(n) - 1)x + \sum_{c \in \mathcal{S}} s_c \right) + \left( \sum_{c \in \mathcal{S}} p_c \right) (-x - s_a) + \left( 1 - p_a - \sum_{c \in \mathcal{S}} p_c \right) \frac{-x}{n+1}.$$

Then for  $n \geq 2$ ,

$$\begin{aligned} & E[a] - E[a : \text{remove } c_1 \text{ from } \mathcal{S}] \\ &= p_a (\gamma(n) - \gamma(n-1)) + p_a s_{c_1} + p_{c_1} (-x - s_a) + \left( 1 - p_a - \sum_{c \in \mathcal{S}} p_c \right) \frac{-x}{n+1} - \left( 1 - p_a - \sum_{c \neq c_1 \in \mathcal{S}} p_c \right) \frac{-x}{n}. \end{aligned}$$

As  $p_{c_1} \geq 0$ ,

$$\left( 1 - p_a - \sum_{c \in \mathcal{S}} p_c \right) \frac{-x}{n+1} - \left( 1 - p_a - \sum_{c \neq c_1 \in \mathcal{S}} p_c \right) \frac{-x}{n}$$

is clearly non-negative. However,

$$p_a s_{c_1} + p_{c_1} (-x - s_a) \geq E[c \text{ against } c_1] \geq 0$$

by assumption. Furthermore,  $\gamma$  is increasing by assumption. One can complete the argument by induction with the  $n = 1$  case holding as

$$E[a \text{ versus singleton}] - E[a \text{ if opposition unfunded}] = E[a \text{ versus singleton}] \geq 0$$

by assumption.

For the second claim, we see

$$E[a] = \sum_{S_j} \frac{\text{Amount funded for } S_j}{\gamma(\#S_j)x + \sum_{c \in S_j} s_c} \left( p_a \left( (\gamma(\#S_j) - 1)x + \sum_{c \in S_j} s_c \right) + \left( \sum_{c \in S_j} p_c \right) (-x - s_a) + \left( 1 - p_a - \sum_{c \in S_j} p_c \right) \frac{-x}{\#S_j + 1} \right).$$

However, we have

$$p_a \left( (\gamma(\#S_j) - 1)x + \sum_{c \in S_j} s_c \right) + \left( \sum_{c \in S_j} p_c \right) (-x - s_a) + \left( 1 - p_a - \sum_{c \in S_j} p_c \right) \frac{-x}{\#S_j + 1} \geq 0$$

by assumption for all  $S_j$ , so  $E[a]$  is a sum of non-negative quantities. □

So, in the third and fourth crowdfunding models above, our expected behaviour is for funders on the side that must be funded as a single option to consider that option they evaluate as being the most likely to win. They will then expect that other funders will not fund very similar options or “clones”<sup>34</sup>. Hence, the funder will tend to either win by default or to be pitted against choice(s) which are meaningfully different from the option he funded. Then Proposition 3 shows us that if a funder Frederick thinks an option is worth insuring against others individually, it will still be worth insuring even if he is pitted against several options collectively. On the other hand, following Proposition 2, funders who fund only a part of the fees for the side of the collective may reasonably fund one or multiple options as they hedge balancing expected returns with variability.

#### 4.8.2.2 Tradeoffs Between Crowdfunding Models

Considering the properties that we showed in Section 4.8.2.1, we summarize some of the various advantages and disadvantages of the four crowdfunding models presented above.

---

<sup>34</sup>Funding such a clone would have a negative expected return, see Proposition 7 in Section 4.11 for analysis of this subject. However, an attacker may be willing to accept a negative expected return in exchange for being able to “grief” the honest funders by also reducing their expected returns. Note, as we see in Proposition 7, if the number of possible clones is small this grief is not very effective.

Model	Advantages	Disadvantages
Require Two Options Funded Individually, Don't Eliminate Other Options	<ul style="list-style-type: none"> <li>• Simple</li> </ul>	<ul style="list-style-type: none"> <li>• More vulnerable to issues where no single option worth funding and/or funding split over multiple options</li> </ul>
Only Funded Options Considered by Jurors in Subsequent Rounds	<ul style="list-style-type: none"> <li>• Less likely to have strategic voting effects in later rounds with fewer options</li> <li>• Better resistance to issues where no single option worth funding</li> </ul>	<ul style="list-style-type: none"> <li>• No natural way to stake on multiple options, so worse effects from liquidity of honest funders being split over multiple options</li> <li>• Pulls system towards Plurality as jurors have to consider which options likely to be funded</li> <li>• Increased risk of nonsense questions being posed to jurors if no good options funded in early rounds</li> </ul>
One option ( $a$ ) must be funded entirely, then funders of other side stake on a set of options $S$	<ul style="list-style-type: none"> <li>• Funders (of options other than <math>a</math>) can make more expressive stakes, deals with liquidity issues</li> </ul>	<ul style="list-style-type: none"> <li>• Asymmetry between funding on two sides of dispute, one side of which must be funded on a single option, may cause UX confusion</li> <li>• Complicates UX to have crowdfunders express which different sets of options willing to fund</li> </ul>
One options ( $a$ ) must be funded entirely, then funders of other side each stake on a set of options $S_j$ , not all necessarily the same	<ul style="list-style-type: none"> <li>• Funders (of options other than <math>a</math>) can make more expressive stakes, deals with liquidity issues</li> <li>• Asking each crowdfunder to specify single set <math>S_j</math> leads to simpler UX than preceding model</li> </ul>	<ul style="list-style-type: none"> <li>• Asymmetry between funding on two sides of dispute, one side of which must be funded on a single option, may cause UX confusion</li> </ul>

### 4.8.3 Parameterization of Stakes

In Section 4.8.2, we saw that funders' returns for financing a successful appeal depend on the arbitration fees that must be paid to the jurors as well as the amounts of additional stake paid by the funders of each side. The choice of stakes in a given application depends on the arbitrable contract being used. In this section, we make some observations that are useful when choosing these values. We generally take the stake required for all options other than the option that won the previous round to be the same. Then funding the option that won the previous round might require less stake.

In the first crowdsourcing model, where only two options can be funded, funding the option  $b$  that won the previous round has a positive expected value if:

$$E[b] = p_b s_a + p_a (-x - s_b) + (1 - p_a - p_b) \frac{x}{2} > 0$$

for any option  $a \neq b$ . If we want to choose the stakes so that it is always worthwhile to fund a previous round winner that is also estimated by a funder to be the option with the single best chances of winning, then one can assume that  $p_b \geq \frac{1}{n}$ , where  $n$  is the total number of possible outcomes. Then the above inequality is satisfied if

$$s_a > s_b + \frac{nx}{2}.$$

Namely, the stake required of the losers of the previous round should be chosen in accordance with the stake required of the previous round winner<sup>35</sup>. Then, as one similarly has:

$$E[a] = p_a s_b + p_b(-x - s_a) + (1 - p_a - p_b) \frac{x}{2} > 0 \Leftrightarrow p_a > \frac{p_b \left( \frac{x}{2} + s_a \right) + \frac{x}{2}}{s_b + \frac{x}{2}},$$

one can choose  $s_a$  so that non-previous round winners are insurable at acceptable thresholds of  $p_a$ .

Reasoning for the other crowdfunding models is similar. In the third and fourth models, it is possible that only one option will be funded in a  $S$ . Based on Propositions 2 and 3, this can be viewed as a worst case scenario for incentivizing fee funders. Hence, it is natural to require the same stakes in this case as in the first model. For the second model, which eliminates outcomes that are not funded, one might want funding the option  $b$  that won the previous round to have a positive expected return as long as  $p_b \geq \frac{1}{\#S}$  for any possible set  $S$  of options that are funded. Then, as all options other than  $b$  require the same stake  $s$ , it is sufficient to have:

$$E[b] = p_b s \#S + (1 - p_b)(-x - s_b) > 0,$$

for which we see that it is sufficient to have  $s > x + s_b$ .

## 4.9 Preventing Pre-Revelation of Juror Votes

### 4.9.1 Commit and Reveal

In Section 4.7.1, we indicated that jurors “commit” to their vote during the voting period. Namely, they submit  $\text{hash}(\text{vote}, \text{salt}, \text{address})$ , and then after the voting period is over, there is a reveal phase where jurors publish their vote and their salt. In this section, we describe this process in greater detail as part of a broader discussion of mechanisms to limit the degree to which information about jurors’ votes circulates during the voting period.

The salt is a random value generated locally in order to add entropy to prevent the use of rainbow tables<sup>36</sup>. During the voting period, jurors can repeat the commitment process to resubmit a new committed vote. This vote might be the same as that of a previous commitment with a different salt, or it may be a new vote. The last commitment received before the end of the commit period will be used for vote aggregation and incentivization, with earlier votes discarded. However, assuming standard collision-resistance properties on the integrity of the hash function, the juror cannot produce distinct vote, salt pairs that produce the same hash. Hence, the juror will not be able to change her vote after the end of the commit period. Through this process, the vote is not visible to other jurors or to the parties. This prevents the vote of a juror from influencing the votes of others.

---

<sup>35</sup>Note the linear dependence of this inequality on the number of possible outcomes. This may be inappropriate for a large number of outcomes, in which case one might choose the stakes so that a previous round winner only is guaranteed to have a positive expected return if it has a larger probability of winning. For example, if one assumes that  $p_b \geq 1/2$ , then  $s_a > s_b + x$  is sufficient. Alternatively, one could slightly modify this model so that in the event that neither  $a$  nor  $b$  wins, the funders of the previous round winner receive back their full deposits with the additional burden of arbitration fees going to the funders of other options. This has the advantage that the relationship between the stakes no longer depends on  $n$ , at the expense of making funding previous round losers even less advantageous and increasing added code complexity.

<sup>36</sup>As currently implemented in the existing Kleros interface, the salt is generated by having the user use their Ethereum private key to sign a block of text that includes an identifier for the specific dispute. Then this salt is always recoverable from the main key by resigning the same text. Nonetheless, the signatures on this text with different keys will be different and are computationally infeasible to obtain without the user’s private key under the assumption of the security of the signature algorithm. Thus, a juror is prevented from copying the commitment of others. Nevertheless, if a user deletes her salt from her local memory, she can regenerate it as long as she retains access to her private key.

## 4.9.2 Penalizing Jurors who Reveal their Vote Too Early

Above we described a commit and reveal scheme that allows jurors to keep their vote hidden until the votes of all jurors are committed to. However, this scheme does not in itself prevent jurors from nonetheless publishing their votes in an attempt to influence other jurors. On some level, one might not expect seeing the vote of a given other juror as being particularly convincing, as even if a given outcome has a publicly visible majority of the votes in the current round, if the outcome is dishonest, jurors who would vote for the minority might expect there to be an appeal. By the “lone voice of reason” effect that we observed in Remark 1, in fact knowing that other voters in your round have voted “incorrectly” gives you *more* of an incentive to vote honestly.

We discuss how to further disincentivize jurors from providing information about how they vote. Note that providing credible information about one’s vote can take a variety of forms:

- Jurors can publish the salt corresponding to their commitment
- Jurors can issue a zero-knowledge proof against their vote commitment that they voted in a certain way
- Jurors can make a smart contract with themselves committing to vote in a certain manner which burns a deposit if they vote differently.

Due to this diversity of strategies jurors can use to pre-reveal their votes, it is impractical to have a totally automated approach to penalize them for doing so (for example, such as a mechanism that would slash a juror’s PNK stake if the salt of her vote is provided prior to the end of the voting period). Instead, in future versions of Kleros, pre-revelation of votes will be an action for which jurors can be challenged in the Process Court as described in Section 4.7.6.

Other approaches that we continue to research for inclusion in future versions of Kleros include those based off the work on anti-pre-revelation games [16]<sup>37</sup>, as well as drawing on the idea of collusion-resistance [19]<sup>38</sup>.

## 4.9.3 Automation of the Revelation of Votes

As this two step processes of committing and then revealing one’s vote described in Section 4.9.1 requires additional user interactions, in some low stakes courts, one might want votes to be issued publicly to simplify the user experience. As described in Section 4.9.2, the fact that Kleros uses an appeal system means that, even if a majority of votes in a voting round have already been cast for a given choice, voting for that choice does not guarantee coherence with the ultimate result used for token redistribution. This limits the effectiveness of a vote copying strategy, and public votes might be acceptable in some cases. Which system is used is determined via a court parameter, see Section 4.12 on governance.

---

<sup>37</sup>In anti-pre-revelation games as they are considered in [16], users can place bets on the votes of jurors and the amount that they are rewarded from a juror deposit or penalized in compensation to the juror depend on the percentage of the voters who vote for each option. This is calibrated in such a way that users can only receive a positive expected return by betting on votes of jurors to the degree that they have some knowledge about how that individual juror will vote *beyond* the broader trends of which options receive more votes on average.

<sup>38</sup>In a collusion resistant mechanism, jurors should not be able to credibly prove to others that they voted in a given way, ideally even after the case is resolved. This property would also have the effect of inhibiting bribe attacks. Note that the ability of jurors in future versions of Kleros to recommit to different votes during the voting period, with only the last vote being counted, is in the spirit if collusion-resistance [19]. Nonetheless, this does not provide complete collusion resistance as jurors can still prove how they voted after the voting period is over. Indeed, as jurors receive different payouts as a function of their vote, it is challenging to create a system where one *cannot* provide evidence that she voted a certain way after the conclusion of a case and distribution of payoffs.

Abstractly, one might hope to have a mechanism where votes just “reveal themselves” at an appointed time without further user involvement. Unfortunately, in a decentralized system where there is no trusted third party who keeps the votes and facilitates this process, this is challenging. Some possible approaches are:

- To have a browser extension or application on the juror’s device that saves the vote and salt information, and then releases it at the appropriate time. Here no one other than the user (and her devices) has access to the vote prior to the reveal period, but one nonetheless avoids the UX issues of requiring the juror to actively issue a second transaction herself. Note, however, that under this approach a user would need to be careful to have her devices turned on at some point during the reveal period.
- To use threshold encryption [48]. Here the juror’s vote and salt can be encrypted under the public key of a group consisting of  $n$  members, any  $t$  of whom can work together to recover the group’s private encryption key and decrypt. See [57] for an implementation of this idea on Ethereum. While this approach imposes some additional overhead, particularly in terms of compensating the participants in the threshold encryption, and raises concerns about the members of the group nonetheless colluding to pre-reveal a vote, the requirements on the actions of jurors in such a system more closely resemble those in a system where votes would just “reveal themselves” without user action.

These solutions might be used in the future to simplify user experience while simultaneously allowing one to more widely enable commit and reveal in different courts. Note that, at least to the degree that third parties addresses other than that of the juror can submit vote, salt information as reveal transactions during the reveal period, from the perspective of the court smart contract it does not matter which of these approaches is used. Hence, regardless of which mechanism is used for vote and salt information to become available during the reveal period, this mechanism could be set up as an overlay over an existing court contract.

## 4.10 Forking

An attacker with a large percentage of active PNK who nevertheless fails a 51% on a given case loses a percentage of the stakes for each time she is drawn. However, as only a subset of PNK are drawn in any given case, even in a late appeal round the attacker would only typically lose a relatively small proportion of her total holdings in lost deposits. In order to maximize the cost of 51% attacks, in this section we will propose an “ultimate vote round” in which all PNK holders vote. This is inspired by a similar mechanism in Augur [49], adapted to the situation of Kleros.

The idea is that one can have a “fork”, creating two versions of the system depending on the result of a particularly contentious vote, where the fork on which an attacker held her tokens would be seen as malicious by the broader market, and her tokens would not be valuable compared to those on the “honest” fork. In current development there is no formal mechanism for facilitating such a fork; however, nothing prevents community members from creating a copy of Kleros in which the result of the contentious decision is reversed and attempt to build a social consensus around using this copy as the main version of Kleros going forward. The remainder of this section details a mechanism that would facilitate the community grouping around forks with desirable properties, as we will see below.

An important feature of Kleros is that cases are often subjective. Hence, it is possible that there will be cases where jurors legitimately disagree. Indeed, we can imagine that a very narrow result in a late appeal round could be the result of any of the following phenomena:

- An attempted 51% attack that is trying to pass through an obviously dishonest result.

- A deep ideological split in how certain types of cases are viewed<sup>39</sup>.
- Honest disagreement on the specifics of a given case that would have little bearing on future cases.

In the first case, one would want there to be a fork to remove the attacker’s influence; in the third case, one would not want there to be a fork as this would unnecessarily fracture the community. A fork over an ideological difference may be justified depending on how serious this difference was viewed by the community. In order to determine what situation we are in, in our “ultimate forking round” PNK holders will specify both their vote in the case as well as information indicating whether they think the case is worth forking over.

Of course, a mechanism that considers a majority vote on whether a case is worth forking over would be ineffective against a successful 51% attack as the attacker would control the outcome of both votes. Instead, PNK holders specify a percentage of total PNK that would go to a given fork at which point they would also be willing to fork. Then a given PNK holder can set this percentage so that her remaining tokens stay with the main fork regardless of the outcome if she thinks the case is a result of an honest disagreement, or so that she definitely forks away from a successful 51% attack. Alternatively, the PNK holder can set this percentage to some intermediate value so that she forks only if there is sufficient support on her side of an ideological split for the new fork to be viable.

We envisage writing arbitrable contracts in such a way that, with unanimous consent of the concerned parties, an arbitrator for an existing contract can be replaced with a fork of Kleros. Hence, this limits the ability of a successful 51% attacker to hold Kleros users hostage except in the relatively rare situation where the attacker has a direct interest. User interfaces can alert concerned parties to the fact that there has been a fork and over what case this fork was made.

#### 4.10.1 Forking Mechanism

We imagine that a given token holder’s utility for a given case with a potential fork as a function:

$$\text{utility} = \text{fct} \left( \begin{array}{l} \text{case outcome} \\ \text{on the fork} \end{array} , \begin{array}{l} \text{percentage of tokens that are on} \\ \text{the same fork as me after the case} \end{array} \right).$$

This allows for a possible trade-off between how egregiously incorrect/unacceptable the winning answer is and the breakdown of how the community splits. We can imagine cases where someone would think that outcome *a* is rather unjust, and it would be worth forking to a universe where outcome *b* won, but only if a large percentage of the community forked with her. Otherwise, if only some marginal amount of the community would have been willing to fork over this case, she prefers tolerating outcome *a* and remaining in the main branch<sup>40</sup>.

We expect that the utility function should be monotonic in the percentage of tokens going to the same fork as you. Namely, all else being equal, we assume that participants would not prefer that the fork they are going to be smaller, as this would be a sign that it would be less likely to catch on. These dynamics are reminiscent of the “battle of the sexes” coordination problem in game theory, where two parties try to coordinate on two possible outcomes, and while they have different preferred outcomes, their preference for landing on the same outcome as the other party is stronger than their preference for their better outcome.

---

<sup>39</sup>For example, the literalist/spirit of the rules type of split that one might have seen with the Augur “who will control the US house after the midterms” market if it had gone on long enough to force a fork [45]. Also compare to the Ethereum/Ethereum Classic fork over how to handle the DAO hack.

<sup>40</sup>This possibly abstracts both the price they might expect the tokens to get on markets going forward after the case as well as their morality/altruism and willingness to participate in a system that they view as just or unjust.



Hence we propose the following: after a case has been appealed the maximum number of times, a final voting round is triggered in which all PNK holder participate<sup>41</sup>. Each PNK holder  $USR_i$  submits  $(a_{ij}, r_{ij}^0) \in L(A) \times [0, 1]^n$ , for  $n \leq \#A$ , that includes a ranked vote  $a_{i1} \geq a_{i2} \geq \dots$  subject to the constraint that  $r_{i1}^0 \leq r_{i2}^0 \leq r_{i3}^0 < \dots$ <sup>42</sup>. The user’s choice of  $r_{ij}^0$  will essentially allow her to specify a minimum threshold for community support for a fork where  $a_j$  is considered the winner at which  $USR_i$  would want to join this fork<sup>43</sup>.

The “main fork”; is the one where the outcome  $a_{\text{main}}$  corresponds to the choice which is selected as the “winner” for settling any payments of ETH for this case and, by default, in other existing contracts. This winning option could be chosen as follows: as participants submit rankings of options  $r_{ij}^0$ , one has enough information about voter preferences to use the same voting system as in Section 4.7.2 to determine a winner in this “ultimate forking vote round”. Then any voter who indicates that  $a_{\text{main}}$  is “acceptable” by including it her ranking will automatically remain on the main fork.

On each fork, the winner on that fork for the purposes of redistribution of PNK for coherence in previous votes is the option selected by that fork. Hence, if a juror in an earlier round believes that a given decision may require a forking round and she has confidence in her ability to choose the fork on which PNK retains market value, the reasoning around the incentivization of the earlier rounds that we have in Section 4.7.3 still holds. Note that it is thus possible that PNK which is “lost” to Bob on a fork where Alice is incoherent may still be held by Alice on the fork where she is coherent. On the other hand, we do not redistribute for the coherence of PNK holders’ “votes” in the forking round except insofar as they determine which fork PNK holders wind up on<sup>44</sup>.

Any PNK that is staked in some court, but for which a set of forking preferences is not provided during the forking period are slashed. Any PNK that is not staked, and for which a forking round vote is not submitted, is sent to the main fork<sup>45</sup>. After the result of each fork is used for PNK redistribution in previous rounds, the tokens that take part in the forking vote receive a bonus such that the total number of tokens on each fork is ultimately equal to the number of tokens on the original fork. Essentially from the point of each fork redistributing the PNK that went to the other forks to the “coherent” PNK that remained. However, unstaked PNK not used to vote do not receive this bonus; this incentivizes all PNK holders to participate in the forking vote. Indeed, it is essential to have very high turnout in any fork vote to maximize resistance to 51% attacks.

Now we will use a procedure that attempts to find the single largest fork(s) compatible with users’

---

<sup>41</sup>Note that, to a great degree, the forking procedure we describe here facilitates what could already be done by the community via social consensus outside of the formal protocol. In the case of social consensus forks it is also possible for a subset of the community to realize that a given outcome is going to win (or already has won), and fork without waiting for and paying appeal fees for a large number of appeals. In future work, we may consider how to expand the forking mechanism to allow a group to coordinate around a fork in an early appeal round, hence facilitating a healthy community in a fuller range of situations where a social consensus fork could have been employed.

<sup>42</sup>This is not optimally expressive; a voter might be willing to fork to  $a_1$  at a lower level of support than she demands to fork to  $a_2$  but nonetheless be more willing to be part of a “large” fork to  $a_2$  than a “small” fork to  $a_1$ , a preference that is not expressible in this model. However, richer expression requires a more complicated user experience and may require a greater running time to resolve the result.

<sup>43</sup>Note that a user specifying a percentage  $r_{ij}^0$  as a threshold for support for a fork at which she is willing to participate is equivalent to specifying a threshold for the number of tokens going to that fork  $r_{ij} \in \mathbb{N}$ ; while it is particularly concrete for users in the interface to think in terms of percentages, we consider these perspectives interchangeably.

<sup>44</sup>As the forking round will not have exactly the same incentive mechanisms as in normal appeal rounds, neither in terms of ETH nor PNK redistribution, one might worry that PNK holders will not give “full rankings”. Indeed, if  $USR_i$  provides a single  $r_{ij}^0 = r_i^0$  for all “acceptable” choices, the system begins to resemble approval voting, which is very sensitive to where jurors draw the line for what is “acceptable” [53]. This may nonetheless be acceptable/appropriate for a forking round vote, as ultimately we want the forking system to be able to allow the community to move on in the case of extremely contentious decisions.





<sup>45</sup>Particularly there is no potential for token holders to wait until they see the market reaction to the different versions of PNK after the vote before making a choice.

expressed preferences<sup>46</sup>. One can perform:

- For all  $i, j$ , compute  $r_{ij} = r_{ij}^0 \cdot \frac{\text{total number}}{\text{of tokens}}$ .
- All token holders  $USR_i$  that rank the option  $a_{\text{main}}$  remain on the fork where  $a_{\text{main}}$  wins.
- For each choice  $a_j \in A$  such that  $a_j \neq a_{\text{main}}$ 
  - Take the list  $L_j$  of voters  $USR_i$  that rank the option  $a_j$  but not  $a_{\text{main}}$ .
  - Take the sum of the number of tokens over all voters on  $L_j$ . Denote this by  $s_j$ .
- For each option  $j$ :
  - Sort  $L_j$  by the user's  $r_{ij}$ .
  - Take the user  $USR_i$  with the largest  $r_{ij}$ . Compare  $r_{ij}$  to  $s_j$ . If  $r_{ij} > s_j$ , remove  $USR_i$  from  $L_j$ .
  - Recalculate  $s_j$  as the sum of the number of tokens over all voters on  $L_j$ .
  - Repeat the two previous steps until  $L_j$  stabilizes.
- Then, over the possible choices  $j$ , take the option that maximizes  $s_j$  and for which  $L_j \neq \emptyset$ .
- All remaining token holders on  $L_j$  go to an alternative fork where  $a_j$  wins.
- Remove the token holders on  $L_j$  and the option  $a_j$  from the originally provided preferences and repeat this process until there is no  $j$  produced such that  $L_j \neq \emptyset$ . This determines if there are still compatible forks possible.

---

<sup>46</sup>Note that in the context of non-binary decisions, there are various ways to translate this information on preferences into a set of consistent forks. Instead of using a rule based on the single largest fork, one could have, for example, attempted to maximize the total number of tokens that fork. This choice has its tradeoffs; for example in maximizing the single largest fork one may have situations where there was an alternative fork that was almost as large that also allowed secondary forks that satisfied a larger percentage of the token pool, an outcome that one might expect to produce higher utility. By choosing this forking rule, one makes a choice that is more easily computable, avoids fragmenting the community more than what is necessary, and has good properties with respect to the presence of clones as we will see below.

	% of Total PNK	Vote	Threshold for willingness to fork to b	Destination Fork
	26%	a	NA	a
	20%	a	NA	a
	12%	b	12%	b
	11%	b	30%	b
	9%	b	40%	a
	9%	a	NA	a
	7%	b	20%	b
	6%	a	NA	a

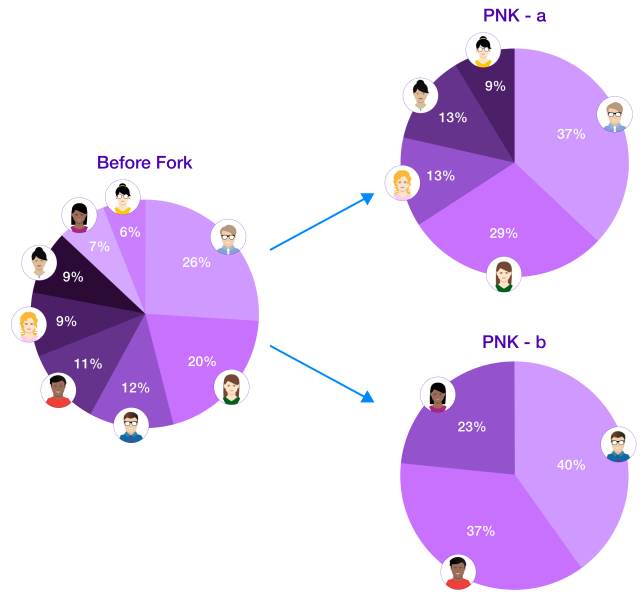


Figure 14: Here a pool of PNK holders vote in a forking round, also expressing their willingness to go to a minority fork. The PNK holders who vote for  $a$  would also generally express a threshold for their willingness to go to a minority fork where  $a$  wins, but as  $a$  already wins the vote here that information is not used and we suppress it for simplicity in the image. Note that one PNK holder voted for  $b$  in the forking round, but as a minority fork to  $b$  involving her would only have 39% support, lower than her threshold, she ultimately remains on the (main)  $a$  fork.

#### 4.10.2 Properties of Forking Proposal

**Proposition 4.** *The above process determines the single largest possible fork.*

*Proof.* We see that the first two loops allow us to determine the largest possible percentage of tokens that can fork to each option  $j$ . Indeed, a simple inductive argument shows that none of the voters removed in the second for loop could have gone to this fork under any subdivision. On the other hand, all voters that remain on  $L_j$  after this process have  $r_{ij} \leq s_j$  as the list is sorted; hence if they all go to the alternative fork together this is compatible with their choices. Then as we maximize over all outcomes, we find the largest possible fork. □

The first for loop takes  $O(\#AN)$  steps, where  $N$  is the number of voters. Sorting  $L_j$  takes  $O(N \log N)$  time. The steps of the second for loop after the sorting require  $O(N)$  time, so the entire for loop takes  $O(N \log N)$  time, repeated  $\#A$  times, so the second for loop takes  $O(\#AN \log N)$  time. The remaining steps require  $O(\#AN)$  steps, and this whole process is repeated at most  $\#A$  times for a total running time of  $O(\#A^2 N \log N)$ .

The above process could be coded directly in the smart contract controlling the forking process. However, to reduce gas costs, it is also possible to use an “optimistic” process where different fork possibilities are submitted during some interval of time, and the contract then determines which produces the single largest fork. Note that a smart contract, given a possible set of forks presented as a list of token holders sent to each outcome with the outcomes pre-sorted by the size of their forks, can determine if it is valid (with respect to users’ expressed preferences) and if it has a greater largest fork than a current choice.

- For each choice  $a_j \in A$  such that  $a_j \neq a_{\text{main}}$

- Take the list  $L_j$  of voters  $USR_i$  who fork to the choice  $a_j$  under the proposed set of forks.
  - Take the sum of the number of tokens over all voters on  $L_j$ . Denote this by  $s_j$ .
  - For each voter  $USR_i$  verify that  $USR_i$  does not rank  $a_{\text{main}}$  and that  $r_{i,a_j} \leq s_j$ . If this is not the case for all  $USR_i$ ,  $a_j$ , return “false”.
- Verify that the  $s_j$  are monotonic, i.e. that the forks were submitted in order by size. (If the  $s_j$  are not monotonic, return “false”.)
  - Starting with the largest  $s_j$ , check that  $s_j \geq f_j$ , the size of the largest fork in the current fork choice. (If the algorithm arrives at a value for which  $s_j < f_j$ , return “false”.) Halt at the first value of  $j$  for which  $s_j > f_j$  and return “true”.
  - If the algorithm does not halt during the comparisons between  $s_j$  and  $f_j$ , i.e. if  $s_j = f_j$  for all  $j$ , return “false”.

The for loop requires  $O(\#AN)$  steps, and then the remaining steps require  $O(\#A)$  steps. Hence, the total running time for this on-chain verification is  $O(\#AN)$  steps.

We introduce the idea of clone independence in the forking process. It is inspired by the idea of clone independence in ranked list voting systems, see [12].

**Definition 1.** *For a given set of preferences expressed by token holders, a set of options  $C = \{a_1, \dots, a_k\}$  is said to be a set of clones if for all token holders  $USR_i$ ,*

- no option outside of  $C$  is ranked between any two options in  $C$ ,
- either all options in  $C$  are ranked or none of them are, and
- $r_{ij_1} = r_{ij_2}$  for all  $a_{j_1}, a_{j_2} \in C$ .

*Then, a forking system is said to be clone independent, given a set of clones  $C$ , deleting the option  $a_j \in C$  from consideration does not change any of the forks produced for options outside of the set  $C$  (either by creating new such forks, by deleting old ones, or by changing which voters are sent to which fork).*

**Proposition 5.** *The single largest fork rule is clone independent.*

*Proof.* Note that for any clone  $a'$  of  $a_{\text{main}}$ , any voter  $USR_i$  who ranks  $a'$ , indicating that it is acceptable, also ranks  $a_{\text{main}}$ , so she is sent to the main fork, which would not be affected by deleting  $a'$ . Then for computing how many tokens are willing to fork to an option in a set of clones  $C = \{a_1, \dots, a_k\}$  such that  $a_{\text{main}} \notin C$ , we see that  $s_{\text{clones}} = s_{j_1} = s_{j_2}$  for all  $a_{j_1}, a_{j_2} \in C$  after the second for loop. Furthermore, deleting one clone does not change  $s_{\text{clones}}$ , nor indeed does it change  $s_j$  for any  $a_j \in A$ . Note, if a fork is created where one of the clones  $a_{j_1}$  wins, then no subsequent option  $a_{j_2} \in C$  can lead to a fork, as the voters on a fork where  $a_{j_2}$  wins would necessarily also have compatible preferences to going to the fork where  $a_{j_1}$  wins. Hence, the order of verification of which lists of voters  $L_j$  go to which forks, other than the at most one fork where an option in  $C$  wins being replaced with a fork to some other option from  $C$ , is unchanged by the deletion of an element of  $C$ . □

**Proposition 6.** *Denote by  $T = \max_{a \in A} \{\# \text{ tokens that rank } a \text{ first}\}$ . Then, under the WoodSIRV voting rule, at least  $T$  tokens go to the main fork.*

*Proof.* Let  $a_t$  be the option that realizes the maximum defining  $T$ . If  $a_t = a_{\text{main}}$ , then in particular at least  $T$  tokens rank  $a_{\text{main}}$ . If  $a_t \neq a_{\text{main}}$ , in the WoodSIRV step that eliminates  $a_t$ , either  $a_{\text{main}}$  must be a Condorcet winner or it must have more first place votes than  $a_t$  after the reallocations of previous steps. In either case, at least  $T$  many tokens must include a ranking for  $a_{\text{main}}$ . □

This proposition can be interpreted in terms the cost to attack the system, particularly how many tokens must be sent to a fork where a malicious outcome has won in order for that outcome to be able to win. We see that WoodSIRV is at least as resistant in this sense as Plurality; however in extreme cases where each voter only ranks a single option and the voting rule degenerates back to Plurality, this result is strict.

## 4.11 Attack Resistance

In order to be a reliable dispute resolution system, Kleros needs to be able to withstand malicious behaviour from participants. In this section, we will discuss the resistance of Kleros to specific attacks that have been identified as being relevant.

### 4.11.1 Buying Half of the Tokens Attack

If a party (or a group of parties colluding) were to buy half of the tokens, it would control the results in the General Court and therefore could ultimately decide all results. However, having a party buying more than half the tokens is highly unlikely if these are fairly distributed.

First, half of the tokens would need to be available for sale, which is not guaranteed. Moreover, the fact that one party could afford all the tokens at current market price does not mean it would be able to buy half of them. Indeed, tokens have increasing marginal costs; they will be dynamically priced on exchanges. Should one party buy a significant percentage of them, the price would go up due to market depth making it increasingly costly to acquire further tokens.

Finally, such an attack could lead to a fork as described in Section 4.10. While this fork would not prevent attackers that successfully obtained half of the tokens from being able to decide the outcome of a given case, it can isolate those attackers' tokens from the rest of the community limiting their resale value and forcing the attacker to absorb their cost. See Figure 15 for a summary of these effects.

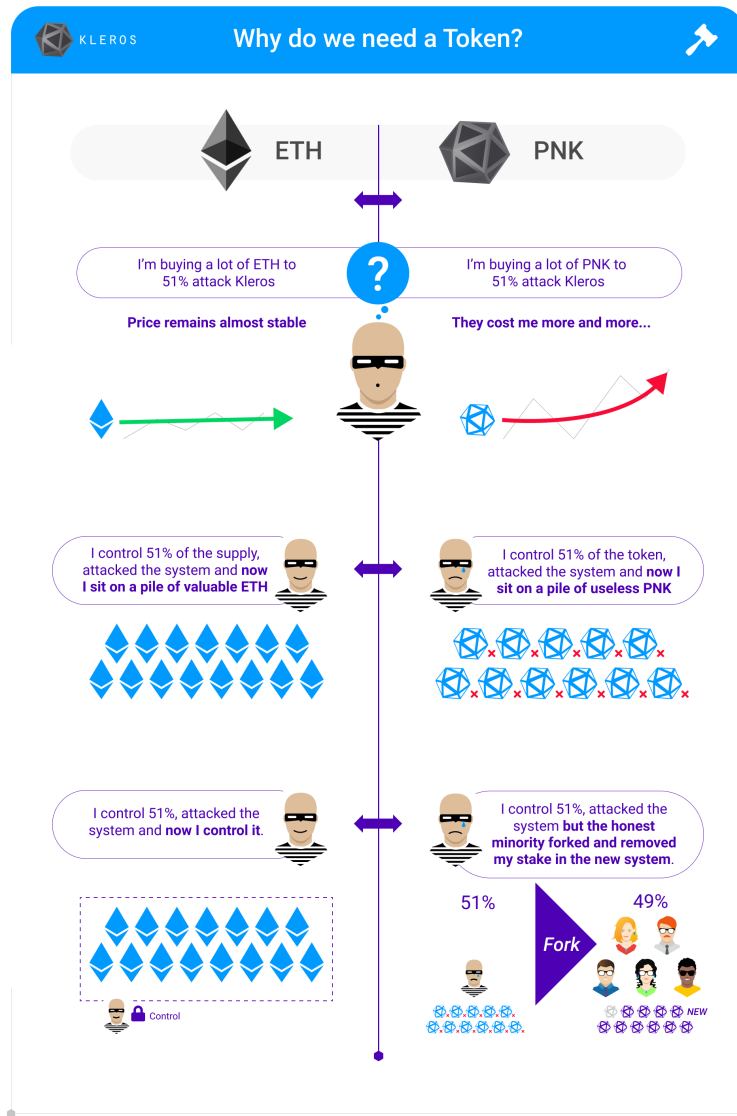


Figure 15: A summary of effects of having a system token, as described in Section 4.5.1, on resistance of Kleros against attacks that attempt to buy up a majority of tokens.

#### 4.11.2 Bribe Attack

Appeals are an important mechanism against bribes. Bribing a small jury is relatively easy. However, since the victim always has the right to appeal, the attacker would have to keep bribing larger and larger juries at a steeply rising cost. The attacker would have to be prepared to spend an enormous amount of money to bribe jurors all the way to the General Court and would very likely lose in the end. To control the verdict of the whole court, the attacker would need to bribe token holders holding more than 50% of the PNK in total<sup>47</sup>.

<sup>47</sup>Note that similarly to the “buying half of the tokens” attack, ultimately the viability of the attack depends on whether 50% of the token pool is corrupted. Here, as the attack does not hold the tokens, the market depth effects discussed above would interact somewhat differently with a potential attack. It is still the case, however, that a successful attack on the General Court would dramatically decrease the value of the tokens (who wants her contracts to be adjudicated by a dishonest court?). Therefore, an attacker should be able to provide more value than 50% of the expected loss from the price drop in order for her bribing offer to be successful (which in almost all cases would exceed

This attack doesn't work in the honest majority model (where more than half of the tokens are controlled by honest parties who won't accept the bribe). But even with a dishonest majority (majority of token holder only searching to optimize their profit), the system can withstand bribing attacks under certain conditions. In practice, a party appealing every decision all the way to the General Court would be extremely unlikely. However, the possibility needs to exist for incentives to be correctly balanced.

### 4.11.3 $p + \epsilon$ Attack

A  $p + \epsilon$  attack is a type of elaborate bribe that promises to pay the bribe only if the attack is unsuccessful, see [17]. Such attacks are particularly designed to target Schelling game systems, warping their incentive structure. These attacks require a high budget but have zero cost if successful<sup>48</sup>. Already in [17] there is a proposed game theoretic response against this attack where jurors use a mixed strategy (jurors only accept the bribe with a defined probability which increases their expected reward compared to accepting the bribe).

Furthermore, we conducted experiments on the Kleros "Doges on Trial" pilot testing user behaviour in the event of a  $p + \epsilon$  attack. See [29] for the results of these experiments. One can also find in [29] comments on how round-based redistribution, which as we saw in Remark 1 leads to the "lone voice of reason" effect, reduces the viability of such attacks.

### 4.11.4 Frivolous Appeals Attack

While the appeal system serves as a defense against bribe attacks, it creates the possibility for attackers with substantial resources to attack the system by appealing cases beyond the point where the other party to case can no longer afford to pay appeal fees<sup>49</sup>. We have detailed in Section 4.8.2 a number of possible mechanisms by which arbitrable contracts can incentivize third party funders to fund the appeals of parties they believe to be honest so that this effect is mitigated.

### 4.11.5 Delay Grief

A related attack would be for an attacker to appeal merely because she wants to delay the execution of a ruling. For example, perhaps she is a competitor of the honest party and hopes to delay when he will receive funds from the case. Unlike in the previous attack, we do not consider the attacker as having the goal of changing the result, so our defenses cannot prevent an attacker from spending money in appeal fees to achieve her desired delay<sup>50</sup>. However, again, as appeal fees increase exponentially over the rounds, an attackers should not be able to maintain this delay for long. Moreover, the General Court has a parameter that sets the maximum number of jurors beyond which further appeals are not possible, see Section 4.8.1. Hence, there is an upper bound on the total amount of time that an attacker can delay a result via this grief.

---

the value at stake in the dispute).

<sup>48</sup>Namely, in case of a successful  $p + \epsilon$  attack, the attacker does not need to make payments to participants and recovers any funds that she locked up. Nonetheless, for the attacker to lock up funds sufficient to convince participants that she is capable of paying all eventual bribes can represent an implicit opportunity cost as the attacker can not use these funds for other purposes during this period.

<sup>49</sup>Note that, as a pure attack, this is only relevant in the appeal fee model where both sides pay appeal fees and the winning party is refunded, see Section 4.8.

<sup>50</sup>An attack where someone accepts to pay a cost to also harm a victim, namely where both the attacker and the victim are worse off, is referred to as a grief, see [18].

### 4.11.6 Clone Funding Grief

We describe a grief on crowdfunding mechanisms, see Section 4.8.2, in the context of non-binary choices. We call this grief “clone funding”. Here the idea is that Frederick funds some (honest) option  $a_j$  to which there is a very similar option  $a_k$ . Funding  $a_k$  is unlikely to be profitable as one would expect that jurors rule for the two options with roughly equal likelihood, however due to the fees that must be paid to the jurors, Frederick and the attacker would then be playing a negative sum game with even likelihood of victory. This grief is relevant in any of the crowdfunding models considered in Section 4.8.2. In the following proposition, for the models that do not allow hedging by funding choices together for the purposes of the following proposition, we take  $\gamma(k)$  to be the number of options that are funded (excluding the one funded by Frederick).

**Proposition 7.** *Assume that all participants possess the same estimates for the probabilities of winning of each outcome. Suppose Frederick funds an option  $a$  that he estimates to have the highest probability of eventually winning. Then the strategy of funding a set of  $k$  other option(s), excluding the previous round winner  $b$  so all options considered require the same stake, with the aim of reducing Frederick’s expected return has a grieving factor of at most  $\frac{1-\gamma(k)+k}{\gamma(k)}$ . In particular, if  $\gamma(k) = \frac{k+1}{2}$ , then these grieving factor are at most one. Moreover, in the fourth model of Section 4.8.2, if  $\gamma(k) = \frac{k+1}{2}$ , then any clonefunding strategy of funding sets of options that do not include  $b$  has a grieving factor of at most one.*

*Proof.* We use the notation of Section 4.8. Suppose Eve funds options  $c_1, \dots, c_k$  each of which has  $p_{c_i} \leq p_a$  and  $s_{c_i} = s_a = s$ . Then, using  $\frac{p_a + \sum_i p_{c_i}}{k+1} \geq \frac{\sum_i p_{c_i}}{k}$  and  $\sum_i p_{c_i} \leq 1 - p_a \leq 1 - \frac{1}{1+k}$ , we have

$$\begin{aligned} E[\text{Frederick}] &= p_a \left( (\gamma(k) - 1)x + \sum_i s_{c_i} \right) + \left( \sum_i p_{c_i} \right) (-x - s_a) - \left( 1 - p_a - \sum_i p_{c_i} \right) \frac{x}{k+1} \\ &\geq s \left( kp_a - \sum_i p_{c_i} \right) + x \left( p_a(\gamma(k) - 1) + \frac{\sum_i p_{c_i}}{k} - \frac{1}{k+1} - \sum_i p_{c_i} \right) \\ &\geq x \left[ \frac{\gamma(k) - 1}{k+1} + \left( 1 - \frac{1}{k+1} \right) \left( \frac{1}{k} - 1 \right) - \frac{1}{k+1} \right] = -x \frac{1 - \gamma(k) + k}{k+1}. \end{aligned}$$

Similarly, using  $p_a \geq \frac{1}{k+1}$ , we have

$$\begin{aligned} E[\text{Eve}] &= \left( \sum_i p_{c_i} \right) s_a + p_a \left( -\gamma(k)x - \sum_i s_{c_i} \right) - \left( 1 - p_a - \sum_i p_{c_i} \right) \frac{kx}{k+1} \\ &= s \left( \sum_i p_{c_i} - kp_a \right) - \gamma(k)x p_a \leq \frac{-\gamma(k)x}{k+1}. \end{aligned}$$

Note these bounds are sharp when  $p_a = \frac{1}{k+1}$ .

Furthermore, under the fourth model of Section 4.8.2 for  $j = 1, \dots, t$ , if Eve provides amounts  $A_j$  of funding to set  $S_j$  such that  $a, b \notin S_j$  and  $p_a \geq p_c$  for all  $c \in S_j$ , then we similarly have

$$\begin{aligned} &E[\text{Frederick}] \\ &= \sum_{S_j} \frac{A_j}{\gamma(\#S_j)x + \sum_{c \in S_j} s_c} \left( p_a \left( (\gamma(\#S_j) - 1)x + \sum_{c \in S_j} s_c \right) + \left( \sum_{c \in S_j} p_c \right) (-x - s_a) + \left( 1 - p_a - \sum_{c \in S_j} p_c \right) \frac{-x}{\#S_j + 1} \right) \\ &\geq \sum_{S_j} \frac{A_j}{\gamma(\#S_j)x + \sum_{c \in S_j} s_c} \left( -x \frac{1 - \gamma(\#S_j) + \#S_j}{\#S_j + 1} \right) \end{aligned}$$



and

$$\begin{aligned}
E[\text{Eve}] &= \sum_{S_j} \frac{A_j}{\gamma(\#S_j)x + \sum_{c \in S_j} s_c} \left( \left( \sum_{c \in S_j} p_c \right) s_a + p_a \left( -\gamma(k)x - \sum_{c \in S_j} s_c \right) - \left( 1 - p_a - \sum_{c \in S_j} p_c \right) \frac{kx}{k+1} \right) \\
&\leq \sum_{S_j} \frac{A_j}{\gamma(\#S_j)x + \sum_{c \in S_j} s_c} \left( \frac{-\gamma(\#S_j)x}{\#S_j + 1} \right).
\end{aligned}$$

However, for  $\gamma(\#S_j) = \frac{\#S_j+1}{2}$ , we have

$$\frac{1 - \gamma(\#S_j) + \#S_j}{\#S_j + 1} = \frac{\gamma(\#S_j)}{\#S_j + 1} = \frac{1}{2},$$

so this also produces a grieving factor of at most one. □

In particular, we note that in the model where only two choices can be financed, clone funding has a grieving factor of at most 1. This can also be obtained in the model that allows hedging by funding options together at the expense of a more complicated function for  $\gamma(k)$ .

## 4.12 Governance Mechanism

As the Kleros protocol gains users and use cases, it will be necessary to create new courts, to make changes in court policies and parameters and to update the platform to new versions with additional features. Such decisions will be made by token holders who have a number of votes equal to the amount of PNK they hold. The governance mechanism can be used to:

1. Set policies: Policies are guidelines about how to resolve disputes. They are the equivalent of the laws in traditional justice systems. They determine which party should win a dispute when particular conditions are met. They can be specific to a particular court.
2. Create new courts.
3. Modify parameters in courts such as:
  - (a) Arbitration fees.
  - (b) Time of each court session.
  - (c) Minimum amount of tokens to be staked.
4. Approve (or remove) alternative templates of voting and incentive systems for use in applications where a different choice of the tradeoffs considered in Section 4.7 is appropriate.
5. Change one of the smart contracts Kleros rely on. This allows arbitrary changes. This can be used for improvements or in an emergency if it appears that some elements of Kleros are not working properly<sup>51</sup>.

Kleros governance decisions are executed on-chain as follows: First, a vote of PNK holders is conducted off-chain on Snapshot [35]. Then lists of transactions to be executed that correspond to the voted upon decisions can be submitted to a governance contract along with a deposit. These lists can be challenged; in this case, one takes advantage of the ability of Kleros itself to provide oracles of off-chain data by requiring jurors in the Kleros Technical Court to determine whether the submitted transactions do, in fact, correspond to what was voted upon on Snapshot.

---

<sup>51</sup>Audits and reviews will be made before the code is deployed. However, it can never be guaranteed 100% that there is not a bug (either in the code or incentives) somewhere. Having this fail-safe provides extra security.

## 5 Additional Future Work

Above, we have already addressed several points on which we intend to improve existing aspects of the protocol in future work. In this section we consider a few other planned improvements not previously discussed.

### 5.1 Redistribution of Funds in Rounds Where No Juror is Coherent

As discussed above, arbitration fees and lost PNK deposits are redistributed on a per round basis between voters that are coherent with the final outcome. If no juror in a given round is coherent, these amounts are currently sent to the governor, and can be allocated by the governance process.

Even when there are only three jurors in a round, it will be relatively rare for a round to have no coherent juror. However, in order for Kleros to be a cost effective alternative to very small scale disputes, such as content moderation, see the examples in Section 6, it will sometimes be necessary to start with initial juries of a single voter. In this case, it is a problem that jurors in this first round will never win PNK deposits from other jurors (because there are no other jurors), however they can nevertheless lose their own PNK deposit to the governor, leading to a reduced incentivization. In future work, we are considering mechanisms by which the governor will automatically redistribute amounts that it accumulates due to rounds where no one is coherent back to the jurors drawn in the corresponding court. Hence, this will average out the lost deposits from situations where jurors attempt to rule honestly but nevertheless wind up being incoherent with the final ruling.

### 5.2 Privacy of Contracts

Solving disputes may require parties to disclose privileged information with jurors. In order to prevent outside observers from accessing this information, in the future, the natural language contracts (English or other) and the labels of the jurors voting options will not be publicly released, and in particular they will not be put on the blockchain. When the contract is created, the creator will submit `hash(contract_text, option_list, salt)` (where `contract_text` is the plain English text of the contract, `option_list` the labels of the options which can be voted by jurors and `salt` is a random number to avoid the use of rainbow tables).

The contract creator will send `{contract_text, option_list, salt}` to each party using asymmetric encryption. This way, parties can verify that the submitted hash corresponds to what was sent to them. In case of a dispute, each party can reveal `{contract_text, option_list, salt}` to jurors which can verify that they correspond to the hash submitted. They can do so using asymmetric encryption such that only the jurors receives the text of the contract and of the options. All these steps will be handled by the application users will run while using Kleros.

### 5.3 Liquid Voting in Governance

In the section on the Governance mechanism above we described how token holders can make a number of decisions for the platform. In this section, we describe a future plan to allow token holders who choose not to vote directly to delegate their vote using a liquid voting mechanism [27]. When a user fails to vote, her voting power will be automatically transferred to her delegate. One can see an illustration of the liquid voting mechanism in Figure 16. This delegation can even be structured so that different delegates are used for different types of votes; for example, for votes regarding the parameters of different courts, allowing one to delegate to a variety of subject-specific experts. Note that delegates

do not need to be humans. They can be smart contract implementing arbitrarily complex voting rules (for example voting on updating fees based on market data).

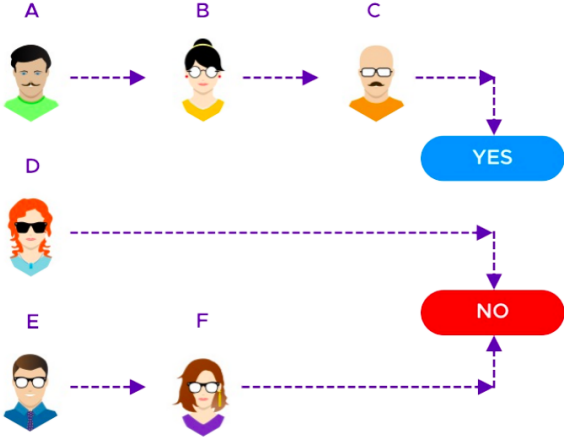


Figure 16: Illustration of a liquid vote

### 5.4 Settlement Logic

There may arise situations where, over the course of a series of appeals, the original parties to a dispute reach an agreement for how they want the original, disputed amount to be divided and wish to prematurely end the dispute resolution process. We intend to add logic to templates of standard arbitrable contracts to allow this<sup>52</sup>.

However, during the course of the dispute other parties, such as jurors and fee funders, will have become financially interested in the result being correct. Indeed, Kleros depends on allowing third parties to finance and call appeals to protect jurors in the event of certain attacks. Hence, it is necessary that the dispute be allowed to continue on, even in the absence of the original parties.

## 6 Applications

Kleros is a general, multipurpose system which can be used in a large number of situations. We present some examples of possible use cases:

- **Escrow:** To pay for an off-chain good or service, the funds can be put in a smart contract. After receiving the good or service, the buyer can unlock the funds to the seller. In case of dispute, Kleros can be used to have the smart contract either reimburse the buyer or pay the seller. Such a Kleros based escrow system is already available, see [33].

Escrows can also be more complex. For example for a rental agreement, the renter can be required to pay a deposit. In case the property is damaged and the renter does not agree on a compensation, a dispute can be created by the owner to claim part of the security deposit.

- **Micro tasking:** Decentralized platforms could pay for microtasks (in the manner of the Amazon Mechanical Turk [1]). Taskers would put a security deposit and submit answers to microtasks. The tasks would be replicated. If a task gets different answers, taskers could admit their mistake, this would transfer a part of security deposit to the taskers who performed the

<sup>52</sup>The current version of the Kleros Escrow has offers a limited structure for settlements, however only before the triggering of a dispute.

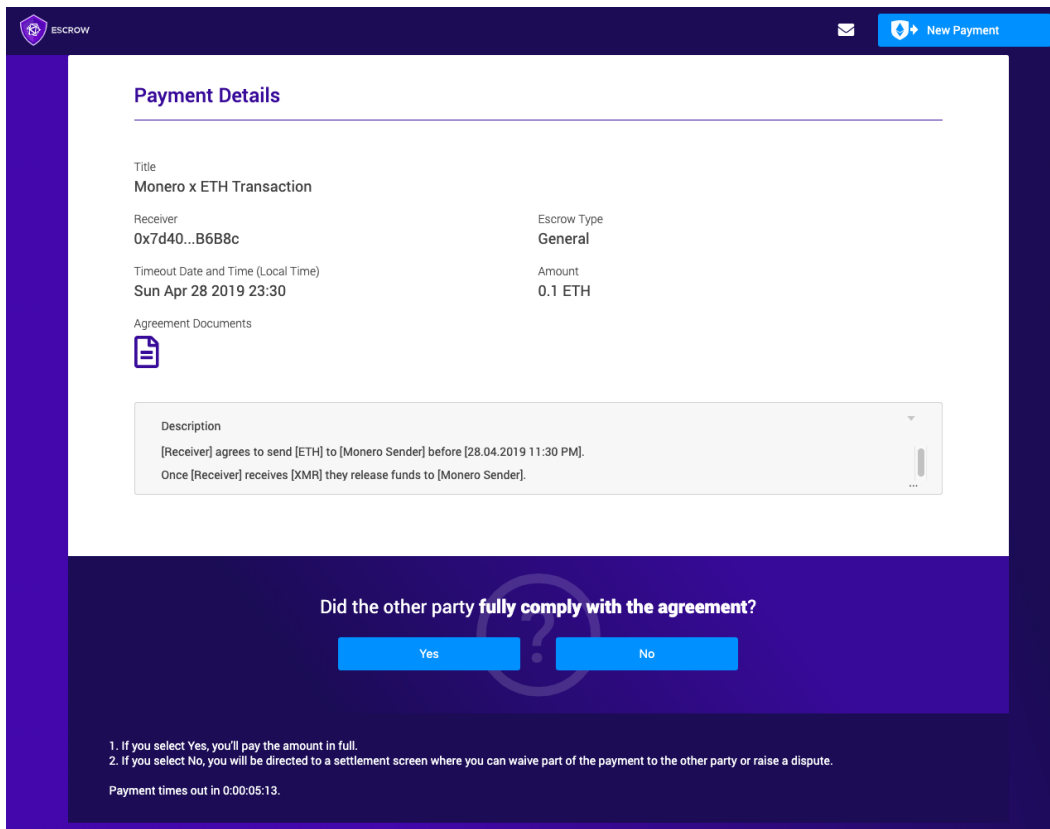


Figure 17: A user of a Kleros escrow deciding whether to raise a dispute.

task correctly. In case multiple taskers stay on their position, a dispute resolution process would ensue and the losing taskers would have part of their security deposit transferred to the winning ones. One example of such a system is Linguo, a decentralized translation platform where the quality of translations can be challenged, resulting in disputes that are judged by Kleros jurors [37].

- **Insurance:** The insuree will pay a fee to the insurer to get a compensation in case a particular event would happen. The insurer will have to put a security deposit which could be common to multiple insurees (respecting risk management rules). When an insured event happens, the insurer can validate it and compensate the insuree. If the insurer does not validate the event, a dispute resolution process would ensue. If the insuree wins the dispute resolution process, funds from the security deposit of the insurer would be transferred to the insuree. In case the security deposit is linked to multiple insurees claiming more than the deposit, a dispute resolution process would also be needed to determine how those funds should be split between insurees.
- **Oracle:** A decentralized data feed to be used by smart contracts was one of the early envisioned use cases of Ethereum [14]. A party (which can be a smart contract) asks a question. Everyone can give a deposit and submit an answer. If everyone gives the same answer, it is returned by the oracle. If there are multiple answers, a dispute resolution procedure ensues. The oracle returns the answer given by the dispute resolution process and parties who put wrong answers lose their deposits which are given to honest submitters. Realitio provides an oracle service that is based on such principles, giving the option to use Kleros for the ensuing disputes [6]. Moreover, other applications that use the Realitio oracle, such as CryptoUnlocked, [44], indirectly depend on this dispute resolution. Particularly, we have researched ways in which such processes can be

efficiently adapted when the oracle is required to output a real-number value, such as in the case of a price oracle [31].

- Curated lists:** Curated lists can be whitelists or blacklists. For example, a whitelist can list smart contracts having undertaken proper audit procedures. A blacklist can list the ENS (Ethereum Name Service [2]) names registered by parties having nothing to do with that name (for example, a malicious party could register “kleros-token-sale.eth”, to scam people into sending funds to that address). Parties could submit items to the list by putting a security deposit. If no one contests that the item belongs to the list for a sufficient amount of time, the name is added and the deposit refunded. If some parties contest by putting a security deposit, a dispute resolution process ensues. If the item is considered belonging to the list, it is added and the submitter gets the deposits of the contesting parties. Otherwise, the deposit of the submitter is given to the contesting parties. Kleros is already being used for a token curated list of tokens that satisfy various properties (for example, such as being ERC20) [34]. One particularly notable application of Kleros is Proof of Humanity, a curated list of distinct individual humans [38], where jurors are required to consider limited biometric information (such as facial features) to mediate challenges that submissions are either of nonexistent people or are of people already on the list. Such a “proof-of-humanity” system allows for a Sybil resistant list of people, which could have applications for quadratic voting schemes [41], more effective airdrops, etc. More generally, the Curate platform [36], provides a registry that acts as a curated list of curated lists. Hence, this system allows individuals to create their own curated lists, subject to some global norms for those lists to be included in a top-level registry.

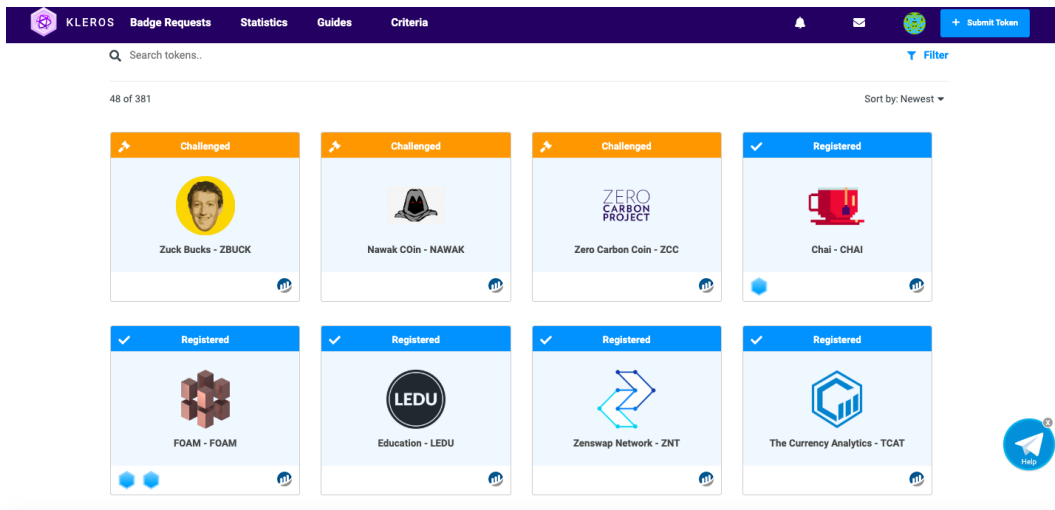


Figure 18: A Kleros-based token curated registry of tokens, where the address and logo of the token must be correct for the token to be allowed on the list.

- Social networks:** Preventing spam, scams and other abuses is a challenge for decentralized social networks. Parties can report violations of the network policies and put a security deposit. If the violation is contested, a dispute resolution process ensues. If it is ruled that no violation happened, the reporter loses her security deposit to the accused party. If the violation is not contested or confirmed by Kleros various effects can be implemented: the content can be removed, the content poster can lose a sign-up deposit and the reach of her other posts can be lowered.

## 7 Conclusion

We have introduced Kleros, a decentralized court system allowing dispute resolution in smart contracts by crowdsourced jurors relying on economics incentives. You can see a summary of how Kleros works in Figure 19.

The rise of the digital economy created labor, capital and product markets that operate in real time across national boundaries. The P2P economy requires a fast, inexpensive, decentralized and reliable arbitration mechanism. Kleros uses game theory and blockchain in a multipurpose dispute resolution protocol capable of supporting a large number of applications in e-commerce, finance, insurance, travel, international trade, consumer protection, intellectual property and academia among many others. Cryptocurrencies are giving many the possibility of having their first bank account to send and receive money in a secure way. Cryptocurrencies are helping millions achieve financial inclusion. Kleros will do the same in access to justice by enabling dispute resolution in a large number of contracts that are too costly to pursue in court. Just as Bitcoin brought “banking for the unbanked”, Kleros has the potential to bring “justice for the unjusted”.

## References

- [1] Amazon mechanical turk. <https://www.mturk.com/>.
- [2] Ethereum name service. <https://ens.domains/>.
- [3] Gnosis. <https://gnosis.pm/>.
- [4] American Bar Association. How courts work. [https://www.americanbar.org/groups/public\\_education/resources/law\\_related\\_education\\_network/how\\_courts\\_work/appeals/](https://www.americanbar.org/groups/public_education/resources/law_related_education_network/how_courts_work/appeals/), 2019.
- [5] Kenneth Arrow. A difficulty in the concept of social welfare. *Journal of Political Economy*, 58, 1950.
- [6] Federico Ast. Kleros-realitio oracle service - getting real information on-chain. Kleros Blog, <https://blog.kleros.io/the-kleros-realit-io-oracle/>, 2019.
- [7] Federico Ast and Bruno Deffains. When online dispute resolution meets blockchain: The birth of decentralized justice. *Stanford Journal of Blockchain Law and Policy*, 2021.
- [8] Manuel Blum. Coin flipping by telephone a protocol for solving impossible problems. *SIGACT News*, 15(1):23–27, 1983.
- [9] Dan Boneh, Joseph Bonneau, Benedikt Bünz, and Ben Fisch. Verifiable delay functions. In *Advances in Cryptology – CRYPTO 2018 - 38th Annual International Cryptology Conference, 2018, Proceedings*, Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), pages 757–788. Springer Verlag, 2018.
- [10] Dan Boneh, Manu Drijvers, and Gregory Neven. *Compact Multi-signatures for Smaller Blockchains: 24th International Conference on the Theory and Application of Cryptology and Information Security, Brisbane, QLD, Australia, December 2–6, 2018, Proceedings, Part II*, pages 435–464. 2018.
- [11] Dan Boneh, Ben Lynn, and Hovav Shacham. Short signatures from the Weil pairing. *Journal of Cryptology*, 17(4):297–319, 2004.

- [12] Felix Brandt, Vincent Conitzer, Ulle Endriss, Jérôme Lang, and Ariel D. Procaccia. *Handbook of Computational Social Choice*. Cambridge University Press, New York, NY, USA, 1st edition, 2016.
- [13] Gilles Brassard, David Chaum, and Claude Crépeau. Minimum disclosure proofs of knowledge. *J. Comput. Syst. Sci.*, 37(2):156–189, 1988.
- [14] Vitalik Buterin. Ethereum, a next-generation smart contract and decentralized application platform. <https://github.com/ethereum/wiki/wiki/White-Paper>, 2014.
- [15] Vitalik Buterin. Schellingcoin: A minimal-trust universal data feed. Ethereum Blog, <https://blog.ethereum.org/2014/03/28/schellingcoin-a-minimal-trust-universal-data-feed/>, 2014.
- [16] Vitalik Buterin. On anti-pre-revelation games. Ethereum Blog, <https://blog.ethereum.org/2015/08/28/on-anti-pre-revelation-games/>, 2015.
- [17] Vitalik Buterin. The  $p + \epsilon$  attack. Ethereum Blog, <https://blog.ethereum.org/2015/01/28/p-epsilon-attack/>, 2015.
- [18] Vitalik Buterin. The triangle of harm. [https://vitalik.ca/general/2017/07/16/triangle\\_of\\_harm.html](https://vitalik.ca/general/2017/07/16/triangle_of_harm.html), 2017.
- [19] Vitalik Buterin. Minimal anti-collusion infrastructure. <https://ethresear.ch/t/minimal-anti-collusion-infrastructure/5413>, 2019.
- [20] Benedikt Bünz, Steven Goldfeder, and Joseph Bonneau. Proofs-of-delay and randomness beacons in Ethereum. [http://www.jbonneau.com/doc/BGB17-IEEESEB-proof\\_of\\_delay\\_ethereum.pdf](http://www.jbonneau.com/doc/BGB17-IEEESEB-proof_of_delay_ethereum.pdf), 2017.
- [21] Craig Calcaterra, Wulf A. Kaal, and Vlad Andrei. Blockchain infrastructure for measuring domain specific reputation in autonomous decentralized and anonymous systems. <https://ssrn.com/abstract=3125822>. U of St. Thomas (Minnesota) Legal Studies Research Paper No. 18-11.
- [22] Lyn Carson and Brian Martin. *Random Selection in Politics*. Praeger, 1999.
- [23] Chainlink Team. Chainlink VRF: On-chain verifiable randomness. <https://blog.chain.link/verifiable-random-functions-vrf-random-number-generation-rng-feature/>, 2020.
- [24] Alisa Cherniaeva, Ilia Shirobokov, and Omer Shlomovits. Homomorphic encryption random beacon. IACR ePrint Archive, <https://eprint.iacr.org/2019/1320.pdf>, 2019.
- [25] John R. Douceur. The Sybil attack. In *Revised Papers from the First International Workshop on Peer-to-Peer Systems*, IPTPS '01, pages 251–260, London, UK, UK, 2002. Springer-Verlag.
- [26] O. Dowlen. *The Political Potential of Sortition: A Study of the Random Selection of Citizens for Public Office*. Luck of the draw : sortition and public policy. Imprint Academic, 2008.
- [27] Bryan Ford. Delegative democracy. <http://www.brynosaurus.com/deleg/deleg.pdf>, 2002.
- [28] David Friedman. A positive account of property rights. *Social Philosophy and Policy*, 11, 1994.
- [29] William George. Doges on trial curated list observations part 2 - deep dive edition. Kleros Blog, <https://blog.kleros.io/cryptoeconomic-deep-dive-doges-on-trial/>, 2018.

- [30] William George. Voting systems for multiple choice Schelling games. <https://github.com/kleros/research-docs/blob/master/multiplechoiceschelling/multiplechoiceschelling3.pdf>, 2020.
- [31] William George and Clément Lesaege. A Smart Contract Oracle for Approximating Real-World, Real Number Values. In *International Conference on Blockchain Economics, Security and Protocols (Tokenomics 2019)*, volume 71 of *OpenAccess Series in Informatics (OASICs)*, pages 6:1–6:15, Dagstuhl, Germany, 2019. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- [32] Allan Gibbard. Manipulation of voting schemes: A general result. *Econometrica*, 41(4):587–601, 1973.
- [33] Stuart James. Kleros Escrow explainer - secure your blockchain transactions today. Kleros Blog, <https://blog.kleros.io/kleros-escrow-secure-your-blockchain-transactions-today/>, 2019.
- [34] Stuart James. Kleros TCR - a deep dive explainer. Kleros Blog, <https://blog.kleros.io/kleros-ethfinex-tcr-an-explainer/>, 2019.
- [35] Stuart James. Handing over the reigns - Kleros governance moves to Snapshot. Kleros Blog, <https://blog.kleros.io/handing-over-the-reigns-kleros-governance-moves-to-snapshot/>, 2020.
- [36] Stuart James. Kleros Curate - the explainer. Kleros Blog, <https://blog.kleros.io/kleros-curate-the-explainer/>, 2020.
- [37] Stuart James. Linguo - the first decentralized translation platform. Kleros Blog, <https://blog.kleros.io/linguo-decentralized-translation-platform/>, 2020.
- [38] Stuart James. Proof of Humanity - an explainer. Kleros Blog, <https://blog.kleros.io/proof-of-humanity-an-explainer/>, 2021.
- [39] Kleros Community. General court policy. <https://court.kleros.io/courts>, 2021. Consulted, June 2021.
- [40] Georgios Konstantopoulos. How does Optimism’s rollup really work? Paradigm Research, <https://research.paradigm.xyz/optimism>, 2021.
- [41] Steven P. Lalley and E. Glen Weyl. Quadratic Voting: How Mechanism Design Can Radicalize Democracy. *AEA Papers and Proceedings*, 108:33–37, 2018.
- [42] L. Laudan. *Truth, Error, and Criminal Law: An Essay in Legal Epistemology*. Cambridge Studies in Philosophy and Law. Cambridge University Press, 2006.
- [43] Clément Lesaege. ERC 792: Arbitration standard. <https://github.com/ethereum/EIPs/issues/792>, 2017.
- [44] Patrick Long. Cryptounlocked oracle upgrade. <https://blog.wetrust.io/cryptounlocked-oracle-upgrade-5c8b22e3375b>, 2019.
- [45] P. H. Madore. Augur House elections market: Alleged reporter says Republicans won the market. <https://www.ccn.com/augur-house-elections-market-alleged-reporter-says-republicans-won-the-market/>, 2018.



- [46] Andrew Mao, Ariel D. Procaccia, and Yiling Chen. Better human computation through principled voting. In *AAAI*, 2013.
- [47] Offchain Labs Team. Chainlink oracles now live on the Arbitrum rollup testnet. <https://offchain.medium.com/chainlink-oracles-now-live-on-the-arbitrum-rollup-testnet-59b7e5d9fed6>, 2021.
- [48] Torben Pryds Pedersen. A threshold cryptosystem without a trusted party. In *Advances in Cryptology — EUROCRYPT '91*, pages 522–526, Berlin, Heidelberg, 1991. Springer Berlin Heidelberg.
- [49] Jack Peterson and Joseph Krug. Augur: a decentralized, open-source platform for prediction markets. <http://bravenewcoin.com/assets/Whitepapers/Augur-A-Decentralized-Open-Source-Platform-for-Prediction-Markets.pdf/>, 2015.
- [50] Jimmy Ragosa. Using Kleros arbitration for dapps on xDai. <https://kleros.gitbook.io/docs/integrations/scalability-and-crosschain/xdai>, 2021. Consulted, June 2021.
- [51] Christian Reitwiessner. From smart contracts to courts with not so smart judges. Ethereum Blog, <https://blog.ethereum.org/2016/02/17/smart-contracts-courts-not-smart-judges/>, 2016.
- [52] Rachel Rothwell. The rise of global litigation funding. <https://www.raconteur.net/risk-management/the-rise-of-global-litigation-funding>, 2017.
- [53] Donald G. Saari and Jill van Newenhizen. Is approval voting an ‘unmitigated evil’?: A response to Brams, Fishburn, and Merrill. *Public Choice*, 59(2):133–147, 1988.
- [54] Mark Allen Satterthwaite. Strategy-proofness and Arrow’s conditions: Existence and correspondence theorems for voting procedures and social welfare functions. *Journal of Economic Theory*, 10(2):187–217, 1975.
- [55] T. C. Schelling. *The strategy of conflict*. Oxford University Press, 1960.
- [56] Markus Schulze. A new monotonic, clone-independent, reversal symmetric, and condorcet-consistent single-winner election method. *Social Choice and Welfare*, 36(2):267–303, 2011.
- [57] Shutter Network Team. Introducing Shutter Network — combating frontrunning and malicious MEV using threshold cryptography. <https://shutter.ghost.io/introducing-shutter-network-combating-frontrunning-and-malicious-mev-using-threshold-cryptography/>, 2021.
- [58] P. Stone. *The Luck of the Draw: The Role of Lotteries in Decision Making*. Oxford University Press, 2011.
- [59] Paul Sztorc. Truthcoin, peer-to-peer oracle system and prediction marketplace. <http://www.truthcoin.info/papers/truthcoin-whitepaper.pdf/>, 2015.
- [60] T.N. Tideman. Independence of clones as a criterion for voting rules. *Social Choice and Welfare*, 4, 1987.
- [61] Sam Vitello, Clément Lesaege, and Enrique Piqueras. ERC 1497: Evidence standard. <https://github.com/ethereum/EIPs/issues/1497>, 2018.

## A Proof of Proposition 1

*Proof.* Denote by  $L$  the number of options. Denote by  $F$  the arbitration fees paid to be split between the voters of this round. For a given winning option  $w$  and collection of votes cast by the other voters in the round, denote by

$$K = \sum_{\mathcal{USR}_k \neq \mathcal{USR}} \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j < w}$$

the total number of pairwise votes on which users other than  $\mathcal{USR}$  are coherent with the ultimate winning choice. Similarly, denote by

$$K' = \sum_{\mathcal{USR}_k \neq \mathcal{USR}} \sum_{a_j \neq w} \mathbf{1}_{\mathcal{USR}_k: a_j \geq w}$$

the total number of pairwise votes on which users other than  $\mathcal{USR}$  are incoherent with the ultimate winning choice. Based on our assumptions, the ultimate winning choice, as well as the votes of the other jurors in the current round, can be considered to be fixed with respect to  $\mathcal{USR}$ 's vote. Hence,  $K$  and  $K'$  are also fixed.

Note that the payoff functions considered are such that providing a vote with ties yields a payoff that is no higher than providing a corresponding strict vote that resolves the ties. Hence providing a vote with ties is a (weakly) dominated strategy. Then suppose  $\mathcal{USR}$  places the ultimate winning choice in the  $i$ th position. Then, as  $\mathcal{USR}$ 's vote is strict and in particular she does not place any other options as tied with  $w$ , she is correct regarding  $L - 1 - (i - 1) = L - i$  pairs and incorrect regarding  $i - 1$  pairs. Hence, when  $\beta = 0$ , her net payoff accounting for lost deposits is:

$$\text{payoff}(i) = \left( [K' + (i - 1)] \frac{D}{L} + F \right) \frac{L - i}{K + L - i} - (i - 1) \frac{D}{L}.$$

In particular, her payoff is a function only of  $i$  (which is notably not the case when  $\beta \neq 0$ ). Note that increasing  $i$  causes  $\mathcal{USR}$  to lose an additional deposit of  $\frac{D}{L}$  which is split between her but also other jurors based on their coherence. Hence,  $\mathcal{USR}$  can only recover a portion of this lost deposit through her reward so

$$\text{payoff}(i) \geq \text{payoff}(i + 1).$$

Then, for any given set of votes by the other voters in the voting round, a standard argument shows that

$$\begin{aligned} E[\text{vote } a_1 > a_2 > \dots > a_L] &= \sum_{j=1}^L \text{payoff}(j) \text{prob}(a_j \text{ wins}) \\ &= \sum_{j=1}^{L-1} \text{prob}(a_j \text{ wins}) \text{payoff}(j) - \text{payoff}(L) (1 - \text{prob}(a_L \text{ wins})) + \text{payoff}(L) \\ &= \text{payoff}(L) + \sum_{j=1}^{L-1} \text{prob}(a_j \text{ wins}) (\text{payoff}(j) - \text{payoff}(L)) \\ &= \text{payoff}(L) + \sum_{j=1}^{L-1} \text{prob}(a_j \text{ wins}) \sum_{i=j}^{L-1} [\text{payoff}(i) - \text{payoff}(i + 1)] \\ &= \text{payoff}(L) + \sum_{i=1}^{L-1} \left( [\text{payoff}(i) - \text{payoff}(i + 1)] \sum_{j=1}^i \text{prob}(a_j \text{ wins}) \right). \end{aligned}$$

is maximized by maximizing  $\sum_{j=1}^i \text{prob}(a_j \text{ wins})$  and hence by voting the options  $a_i$  in order by their probability of winning. As this is true for any given set of votes by the other voters, and these votes are independent of  $\mathcal{USR}$ 's vote and the eventual outcome, we have the desired result.  $\square$

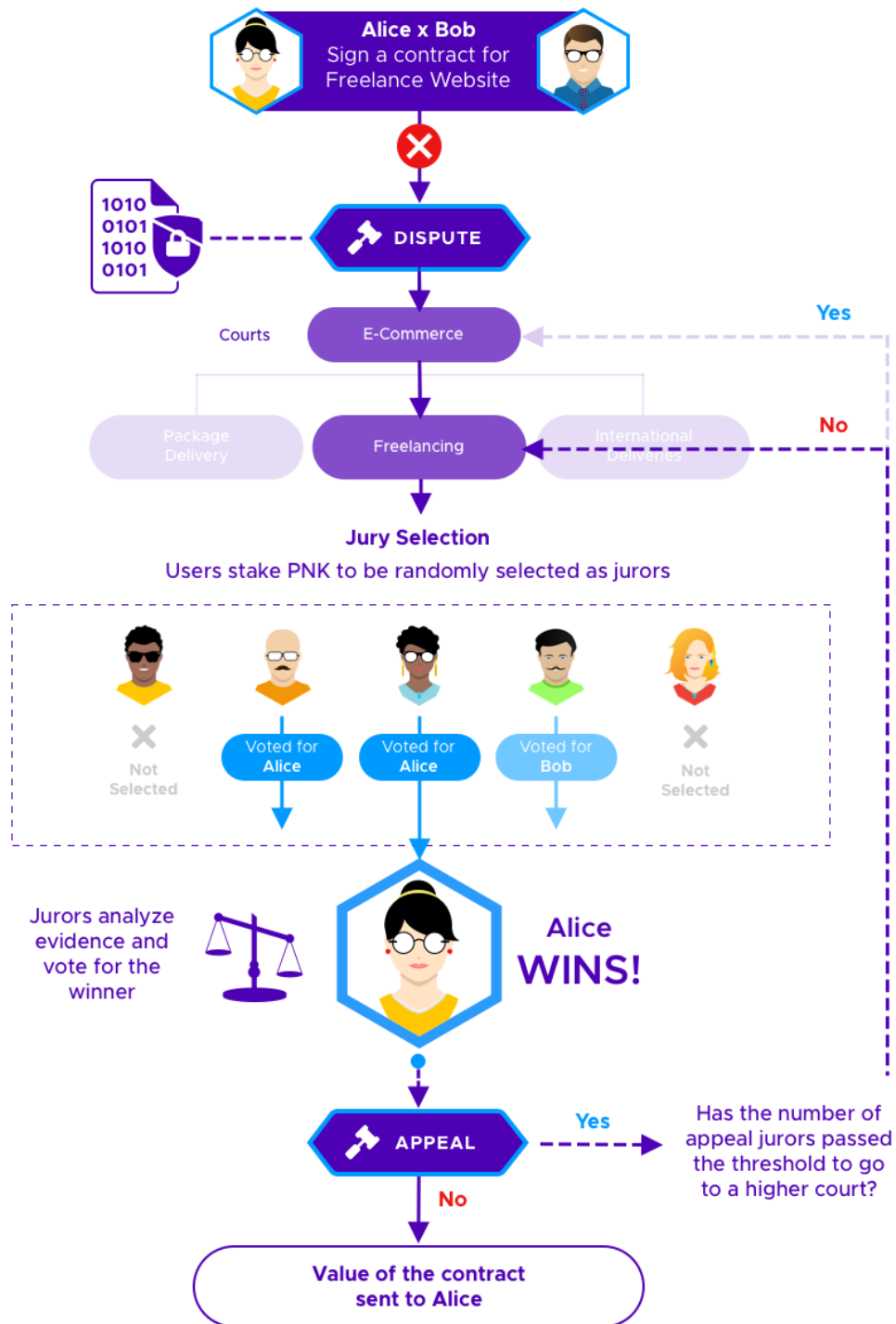


Figure 19: Example of dispute summing up how Kleros works.