

# How to Incentivize Data-Driven Collaboration Among Competing Parties

[Extended Abstract]<sup>\*</sup>

Pablo Daniel Azar  
MIT  
77 Massachusetts Ave  
Cambridge, MA 02139, USA  
azar@csail.mit.edu

Shafi Goldwasser  
MIT\* / Weizmann<sup>†</sup>  
\*77 Massachusetts Ave  
Cambridge, MA 02139, USA  
<sup>†</sup>Herzl St 234  
Rehovot, 7610001, Israel  
shafi@csail.mit.edu

Sunoo Park  
MIT  
77 Massachusetts Ave  
Cambridge, MA 02139, USA  
sunoo@csail.mit.edu

## ABSTRACT

The availability of vast amounts of data is changing how we can make medical discoveries, predict global market trends, save energy, and develop new educational strategies. In certain settings such as Genome Wide Association Studies or deep learning, the sheer size of data (patient files or labeled examples) seems critical to making discoveries. When data is held distributedly by many parties, as often is the case, they must share it to reap its full benefits.

One obstacle to this revolution is the lack of willingness of different entities to share their data, due to reasons such as possible loss of privacy or competitive edge. Whereas cryptographic works address the privacy aspects, they shed no light on individual parties' losses and gains that may influence their decision to collaborate. This is the question we model and study in this paper. The *order* in which collaborators receive the outputs of a collaboration will be a crucial aspect. Our contributions are:

- We formalize a model of  $n$ -party collaboration for computing functions over private inputs, in which the participants receives their outputs in order (i.e. player  $i$  receives his output in position  $\pi(i)$ , where  $\pi$  is a permutation over  $n$  elements) and where each output “improves” on all previous outputs. Both  $\pi$  and the outputs received depend on all players' input data. The model includes a reward function that captures the utility of collaboration for player  $i$  depending on his input, his output, other players' outputs, and  $\pi(i)$ .
- We design mechanisms to compute a distribution of outputs and an ordering of output delivery, based on the  $n$  participants' private inputs. We say that a col-

laboration outcome is a *collaborative equilibrium* if it guarantees a higher reward for all participants when joining a collaboration compared to not joining it. We show that while in general computing a collaborative equilibrium is NP-complete, we can design polynomial-time algorithms for computing it for a range of natural model settings.

- We show how to extend the theory of multi-party computation (MPC) to impose an ordering on the delivery of outputs to different players. This extension, named *ordered MPC*, can be used to implement our mechanisms without a centralized trusted party and in the presence of a subset of colluding players who may deviate from the protocol in an arbitrary fashion, under cryptographic assumptions. Since the delivery schedule depends on players private inputs, it may leak information. Hence, we enhance the MPC privacy requirement to require each player to learn nothing more than his output and his own position in the output ordering. We also adapt the notions of guaranteed output delivery and fairness to require instead *ordered output delivery* and *prefix-fairness*.
- We introduce *timed-delay MPC*, where explicit time delays are introduced into the output delivery schedule. We use the classical time-lock primitive, to obtain protocols for timed-delay MPCs. Finally, we introduce a new concept called *time-line puzzles*, which are a natural extension of classical time-lock puzzles, in which multiple data items can be locked so that their eventual openings can be “serialized” in time.

<sup>\*</sup>A full version of this paper is available at <http://eprint.iacr.org/2015/178>

# 1. INTRODUCTION

The availability of vast amounts of data is changing how we can make medical discoveries, predict global market trends, save energy, improve our infrastructures, and develop new educational strategies. Indeed, it is becoming clearer that *sample size* may be the most important factor in making surprising new discoveries in a number of areas such as *genome-wide association studies*<sup>1</sup> (GWAS) and *machine learning* (ML), as witnessed by the striking success of GWAS studies with large samples for schizophrenia<sup>2</sup> [BP12; PGC14; For11] and the success of deep learning in ML.

When large data is required, often parts of the data are held by different entities. Such entities need to share their data, or at least engage in a collaborative computation where each entity manages its own private data, in order for society to reap the benefit of large sample sizes. Referring back to the GWAS example, success was explicitly attributed to such collaboration: “The schizophrenia study was made possible due to unusually large scale collaborations among many institutes... This level of cooperation between institutions is absolutely essential... If we are to continue elucidating the biology of psychiatric disease through genomic research, we must continue to work together.” [Ins14]

Unfortunately, the above example is the exception rather than the rule. A major obstacle to the big-data revolution is the lack of willingness of different entities to share data in collaborations with each other: so-called “data hoarding”. One obstacle is privacy concerns, where parties refuse to collaborate, in order to protect the privacy of their data. Privacy, however, is not the only obstacle.

An equally important obstacle is competition between entities holding data. When access to data carries tangible rewards, say, if the entities are companies competing for a share of the same market or research laboratories competing for scientific credit, it is unclear whether an individual collaborator is better off, even if it is clear that better overall conclusions can be drawn from collaboration. Stated in more game-theoretic terms, the entities face the following dilemma: *whereas the overall societal benefit of collaboration is clear, the utility for an individual collaborator may be negative, so why collaborate?* Addressing this question is the topic of this paper.

In this paper, we present a formal model for collaboration in which this question can be analyzed, as well as design mechanisms to enable collaboration where all collaborators are provably “better off”, when possible. The *order* in which collaborators receive the outputs of a collaboration will be a crucial aspect of our model and mechanisms. We believe that timing is an important and primarily unaddressed issue in data-based collaborations. For example, in the scientific research community, data sharing can translate to losing a prior publication date. In financial enterprises, the timing of investments and stock trading can translate to large financial gains or losses.

<sup>1</sup>A genome-wide association study is an investigation of common genetic variants in a population, in order to identify genetic variants that are associated with a given trait.

<sup>2</sup>“Dramatic increase in patient data size enabled the discovery of more than 100 gene loci associated with the disease up from a handful loci seen with small sets of patients. This was made possible due to an unusually large scale collaborations among many institutes.”

The collaboration mechanisms we develop assume a central trusted party; however, we will show that this assumption is not necessary under standard cryptographic assumptions. We show in Section 3 how the mechanisms can be implemented in a decentralized way by  $n$  distrustful parties by using new extensions of classical secure multiparty computation that give formal guarantees on the *order and timing* of output delivery as well as the standard requirements of privacy and correctness.

## 1.1 Summary of our contributions

### 1.1.1 A model of collaboration

We propose a model for collaboration which enables the determination of whether the utility obtained by a collaborator outweighs the utility he may obtain without collaboration. The ultimate desired outcome of a collaboration is to learn a parameter of the (unknown) joint distribution from which the participants’ input data  $x_1, \dots, x_n$  is drawn. This can be expressed as  $y^* = f(\mathcal{X})$  where  $\mathcal{X}$  is the joint distribution of input data and  $f$  is a known function. In our model, the outcome of a collaboration is a pair  $(\pi, \vec{\mathcal{Z}})$  where  $\pi$  is a permutation of player identities and  $\vec{\mathcal{Z}} = (\mathcal{Z}_1, \dots, \mathcal{Z}_n)$  where each  $\mathcal{Z}_{\pi(i)}$  is a distribution that corresponds to player  $i$ ’s “estimate” of  $y^*$ . We think of  $\mathcal{Z}_{\pi(i)}$  as the public output of player  $i$ : for example, in the setting of scientific collaboration,  $\mathcal{Z}_{\pi(i)}$  would be player  $i$ ’s academic publication. Our model setup assumes an underlying score function which assigns scores to the players’ outputs.

The model includes a *reward function*  $R_t$  which characterizes the gain in utility for any given party  $i$  in a collaboration. The reward that a party  $i$  gets depends on how much his score  $s(\mathcal{Z}_{\pi(i)})$  *improves on* the previous state of the art  $s(\mathcal{Z}_{\pi(i)-1})$ , and on  $\pi(i)$ , namely, *when* the party makes his public output. Specifically, the reward function includes a multiplicative *discount factor*  $\beta^t$  where  $\beta \in [0, 1]$  and  $t$  is the time of publication, meaning that the reward from a publication is “discounted” more as time goes on.

$$R_t(\pi, \vec{\mathcal{Z}}) = \beta^t \cdot (s(\mathcal{Z}_{\pi(t-1)}) - s(\mathcal{Z}_{\pi(t)}))$$

To determine whether the utility of collaboration outweighs the utility of working on one’s own, our model uses “outside payoff” values  $\alpha_i$  which are the score that party  $i$  would obtain *without collaborating*.  $\alpha_i$  can be computed directly from the input  $x_i$  of party  $i$ .

### 1.1.2 Mechanisms and collaborative equilibrium

We define a notion of *collaborative equilibrium* in which all parties are guaranteed a non-negative reward, and develop *mechanisms* for collaboration that compute such equilibria (when they exist). When an equilibrium exists, our mechanism delivers a sequence of progressively improving approximations of  $y^*$  to the collaborating parties. More specifically, the mechanism will take as input the data of all parties, and output a pair  $(\pi, \vec{\mathcal{Y}})$  where  $\pi$  is a permutation of player identities and  $\vec{\mathcal{Y}} = (\mathcal{Y}_1, \dots, \mathcal{Y}_n)$  specifies the outcomes to be delivered to the players: each  $\mathcal{Y}_{\pi(i)}$  is the approximation to  $y^*$  that is given to player  $i$  at time-step  $\pi(i)$ , such that the score of the outputs is increasing with time. That is,  $s(\mathcal{Y}_{\pi(1)}) > \dots > s(\mathcal{Y}_{\pi(n)})$ . We emphasize that both the order  $\pi$  and the outputs  $\mathcal{Y}_i$  are computed based on the inputs of all players.

When player  $i$  receives an output  $\mathcal{Y}_{\pi(i)}$  from the central mechanism, she may combine  $\mathcal{Y}_{\pi(i)}$  with the information that she learned from prior public outputs and her own input  $x_i$ , to generate a public output  $\mathcal{Z}_{\pi(i)}$ . We first prove that the ability of the players to learn from others' publications, in general, will make the problem of deciding whether there exists an equilibrium is NP-complete (see Theorem 3).

Next, we show that there is a polynomial-time mechanism that can output an equilibrium whenever one exists (or output **NONE** if one does not exist) for a variety of model settings and parameters which we characterize (see Theorem 2). An example of a setting when a polynomial-time mechanism is possible is when

- an upper bound  $\mu_j$  can be set on the amount of information that any player can learn from a given player  $j$ 's publication, and
- it is possible to efficiently compute, for any  $y^*$  and  $\delta > 0$ , an "approximation"  $\mathcal{Y}'$  such that  $s(\mathcal{Y}') = \delta$ .

In a nutshell, the bounds  $\mu_j$  can be used to define a weighted graph in which the weight of the minimum-weight perfect matching determines the existence of a collaborative equilibrium.

### 1.1.3 Cryptographic protocols to implement the mechanisms

We develop cryptographic protocols for implementing the mechanisms without a centralized trusted party and in the presence of a subset of colluding players who may deviate from the protocol in an arbitrary fashion, under cryptographic assumptions. The protocols compute the collaboration outcome  $(\pi, \mathbf{y})$  via multi-party secure computation on players' private inputs. Since a crucial aspect of the mechanism's ability to yield non-negative reward to all players is the delivery of outputs in order, we need to extend the classical notion of MPC to incorporate guarantees on the order and timing of output delivery. These extensions may be of interest independent of the application of mechanisms for incentivizing collaborations.

We define *ordered MPC* as follows. Let  $f$  be an arbitrary  $n$ -ary function and  $p$  be an  $n$ -ary function that outputs permutation  $[n] \rightarrow [n]$ . An ordered MPC protocol is executed by  $n$  parties, where each party  $i \in [n]$  has a private input  $x_i \in \{0, 1\}^*$ , who wish to securely compute  $f(x_1, \dots, x_n) = (y_1, \dots, y_n)$  where  $y_i$  is the output of party  $i$ . Moreover, the parties are to receive their outputs in a particular *ordering* dictated by  $p(x_1, \dots, x_n) = \pi$  where  $\pi$  is a permutation of the player identities. Since the choice of  $\pi$  depends on private inputs, it may leak information: hence, we formulate an enhanced *privacy* requirement for ordered MPC that each player should learn his output and his *own* position in the output ordering, and nothing more (see Definition 9).

We show a simple transformation from classical MPC protocols for general functionalities  $f$  to ordered MPC protocols for general functionalities  $f$  and permutation functions  $p$  that achieve enhanced privacy, even when a minority of the  $n$  players may be colluding to sabotage the protocol (see Theorem 4). The assumptions necessary are the same as for the classical MPC constructions (e.g. [GMW87]). When the colluding players are in majority, it is well known that output delivery to all honest parties cannot be guaranteed [Cle86].

Next, we define *timed-delay MPC*, where explicit time delays are introduced into the output delivery schedule. Time delays between the outputs may be crucial to enable parties to reap the benefits of their position in the order. We give two constructions of timed-delay MPC in the honest majority setting<sup>3</sup>. First, we give a conceptually simple protocol which runs "dummy rounds" of communication in between issuing outputs to different players, in order to measure time-delays. The simple protocol has the flaw that all (honest) players must continue to interact until the last party receives his output (that is, they must stay online until all the time-delays have elapsed). To address this issue, we present a second protocol assuming the existence of time-lock puzzles [RSW96] in addition to the classical MPC [GMW87] assumptions (see Theorem 7). Informally, a time-lock puzzle is a primitive which allows "locking" of data, such that it will be released after a pre-specified time delay, and no earlier. Our second timed-delay MPC protocol, instead of issuing outputs to players in the clear, gives to each party his output *locked* into a time-lock puzzle; and in order to enforce the desired ordering, the delays required to unlock the puzzles are set to be an increasing sequence. An issue that arises when giving out time-lock puzzles to many parties is that different parties may have different computing power, and hence solve their puzzles at different speeds: for example, it is clear that we cannot guarantee that players learn their outputs in the desired ordering if some players compute arbitrarily faster than others. Still, we show that our protocol is secure and achieves ordered output delivery in the case that the difference between any two players' computing power is known to be bounded by a logarithmic factor. If the assumption about computing power does not hold, then the protocol still achieves security (i.e. correctness and privacy), but the ordering of outputs is not guaranteed.

The definition of ordered and timed-delay MPC inspire new notions unrelated to the central topic of this paper. In particular:

- **Time-lines.** Inspired by the application of time-lock puzzles to time-delayed MPC, we propose the new concept of a *time-line*, where multiple data items can be locked so that their unlocking must be serialized in (future) time. See Section 4.4 for details.
- **Prefix-fairness.** In the traditional MPC landscape, fairness is the one notion that addresses the idea that either all parties participating in an MPC should benefit, or none should. Fairness requires that either all players receive their output, or none do. However, it is well-known that fairness is achievable when a majority of the players are honest, but it is *not* achievable for general functionalities when a majority of players are faulty [Cle86]. We propose a refinement of the classical notion of fairness in the setting of ordered MPC, called *prefix-fairness*, where players are to receive their outputs one after the other according to a given ordering  $\pi$ , and the guarantee is that *either* no players receive an output *or* those who do strictly belong to a prefix of the mandated order  $\pi$  (see Definition 11). Prefix-fairness can be achieved for general functionalities and

<sup>3</sup>We cannot hope to achieve timed-delay MPC in the case of dishonest majority since, as mentioned in the preceding paragraph, even output delivery cannot be guaranteed in this setting.

any number of faulty players, under the same assumptions as classical MPC [GMW87] (see Theorem 5).

## 1.2 Discussion and interpretation of our work

### *Slowing down scientific discovery?*

Intuitively, the mechanisms we develop always take the following form: the mechanism computes the “best possible estimate”  $\mathcal{Y}^*$  of  $y^*$  given the input data of the players, and then hands out a sequence of successively more accurate (according to the score function) outcomes, where the final party receives  $\mathcal{Y}^*$ .

One may ask: why slow down scientific progress and hand out inferior results when better ones are available? We argue that progress will in fact be *enhanced*, not slowed down, by this methodology, as it will be a decisive factor in parties’ willingness to collaborate in the first place. This bears great similarity to the original philosophy of *differential privacy* and privacy-preserving data analysis more generally. In these fields, accuracy (so-called utility) of answers to aggregate queries over items in database is partially sacrificed in order to preserve privacy of individual data items, as a way to encourage individuals to contribute their data items to the database. In an analogous way, in order to get results based on the large data sets held by potential collaborators, we sacrifice the *speed* of discovery of the “ultimate” collaboration outcome: we are willing to pay this price to incentivize parties to collaborate and contribute their data. In contrast to differential privacy, we do not sacrifice ultimate accuracy. The last collaborator to receive an output, receives the ideal outcome  $\mathcal{Y}^*$ . Namely,  $\mathcal{Y}_n = \mathcal{Y}^*$ .

### *The Fort Lauderdale example: the importance of time.*

A recurring idea in this work is the importance of time and ordering of research discoveries, which is inspired in part by the following striking example from the field of genomics. In the 2003 Fort Lauderdale meeting on large-scale biological research [Wel03], the gathering of leading researchers in the field recognized that “pre-publication data release can promote the best interests of [the field of genomics]” but “might conflict with a fundamental scientific incentive – publishing the first analysis of one’s own data”. Researchers at the meeting agreed to adopt a set of principles by which although data is shared upon discovery, researchers hold off publication until the original holder of the data has published a first analysis. Being a close-knit community in which reputation is key, this was a viable agreement which has led to great productivity and advancement of the field. However, more generally, their report states that “incentives should be developed by the scientific community to support the voluntary release of [all sorts of] data prior to publication”. This example teaches us to focus on three key aspects of collaboration: the incentive to collaborate has to be clear to all collaborators; there must be a way to ensure adherence to the rules of collaboration; and timing is of the essence.

### *Privacy implies increased utility.*

Although the goal of our work is to design mechanisms to *incentivize collaboration* by increasing the utility of collaborations rather than focusing on the privacy of individual entities’ input data, MPC protocols prove to be an important technical tool to implement the mechanisms which guarantee

increased utility. As a by-product, the use of MPC provides our mechanisms with the additional guarantee of privacy.

### *Future directions.*

When collaboration is feasible, each party  $i$  in our model is guaranteed a reward from collaborating that is greater than the reward  $\alpha_i$  that they could get on their own. However, the contributions of the players’ data to the computation of the final output  $\mathcal{Y}^*$  (i.e. the best possible estimate of  $y^*$  given the joint input data) may be asymmetric: some special player  $i^*$  may have some data that helps solve the “puzzle”, but this player  $i^*$  may not be known a priori before the participants decide to collaborate<sup>4</sup>. An interesting future direction would be to develop mechanisms where, even if we do not know a priori which players have higher quality data, we can still allocate more credit after the collaboration is done to the players whose contribution turned out to be the most valuable.

Another future direction of interest to design *truthful* mechanisms so that collaborating parties will be provably incentivized to submit their true and accurate data as input. In our work, we assume that, while we can incentivize the players to collaborate or not, once they decide to collaborate they are truthful about the value of their dataset  $x_i$ . From the point of view of scientific publications, this assumption is reasonable if we believe that the experiments that generate this data can be verified or replicated, and that a failure to replicate would hurt a scientific group’s reputation. However, there are many settings, such as businesses pooling their data together to generate larger profits, where the parties may be incentivized to lie about their output  $x_i$ . Since we are already assuming that parties are rational, a future direction would be to develop mechanisms where, even when parties can lie about  $x_i$  (because  $x_i$  cannot be verified by others), they are still incentivized to report it truthfully. One possible direction is where  $x_i$  is the output of some long  $\#P$  computation (for example, a Markov Chain Monte-Carlo simulation), where (a) replicating the computation would take a very long time and delay publication for everyone in the group and (b) player  $i$  cannot prove in a classical way that their output  $x_i$  is correct. Even in this case, player  $i$  can be incentivized to give the right answer via a rational proof [AM12; AM13; GHRV14].

Our setting is most useful and likely to lead to collaboration when the marginal value of data is increasing. Interestingly, this does not include classical learning problems that can be stated as minimizing a convex loss function (or maximizing a concave value function). However, as new technology is developed, many new interesting machine learning problems with non-convex objectives have become practical. Bengio and LeCun [BL+07] give more thorough discussion of this situation.

## 1.3 Other related work

The problem of how to make progress in a scientific community has been studied in other contexts. Banerjee, Goel and Krishnaswamy [BGK14] consider the problem of partial

<sup>4</sup>An example in the same vein is the following. In the medical setting, a hospital with a larger patient population will clearly have more patient data than a small facility, and yet access to data of small but homogeneous or rare communities can at times be more valuable than access to larger heterogeneous sets of data.

progress sharing, where a scientific task is modeled as a directed acyclic graph of subtasks. Their goal is to minimize the time for all tasks to be completed by selfish agents who may not wish to share partial progress.

Kleinberg and Oren [KO11] study a model where researchers have different projects to choose from, and can work on at most one. Each researcher  $i$  has a certain probability of being able to solve a problem  $j$ , and she gets a reward  $w_j$  if she is the only person to solve it. If multiple researchers solve the problem, they study how to split the reward in a socially optimal way. They show that assigning credit asymmetrically can be socially optimal when researchers seek to maximize individual reward, and they suggest implementing a “Matthew Effect”, where researchers who are already credit-rich should be allocated more credit than in an even-split system. Interestingly, this is coherent with the results of our paper, where it is socially optimal to obfuscate data so that researchers who are already “ahead” (in terms of data), end up “ahead” in terms of credit.

Cai, Daskalakis and Papadimitriou [CDP14] study the problem of incentivizing  $n$  players to share data, in order to compute a statistical estimator. Their goal is to minimize the sum of rewards made to the players, as well as the statistical error of their estimator. In contrast, our goal is to give a decentralized mechanism through which players can pool their data, and distribute partial information to themselves in order so as to increase the utility of every collaborating player.

Boneh and Naor [BN00] construct timed commitments that can be “forced open” after a certain time delay, and discuss applications of their timed commitments to achieve fair two-party contract signing (and coin-flipping) under certain timing assumptions including bounded network delay and the [RSW96] assumption about sequentiality of modular exponentiation.

### Roadmap.

Section 2 covers the scientific collaboration model, mechanisms, and feasibility theorems. Section 3 covers the definitions and constructions of ordered MPC, and Section 4 covers definitions and constructions of timed-delay MPC, and associated primitives such as time-line puzzles.

## 2. SCIENTIFIC COLLABORATION MODEL

In this section, we present and analyze mechanisms for scientific collaboration in our model. In our exposition, we focus primarily on the setting of scientific collaboration and publication. However, we want to highlight that our results apply to more broad collaboration and discovery in general, in which case a “publication” should be thought of as any kind of public output.

### Notation.

We denote by  $[n]$  the set  $\{1, \dots, n\}$  of integers between 1 and  $n$ , and by  $[n] \rightarrow [n]$  the set of all permutations of  $[n]$ . For a set  $X$ , we write  $\Delta(X)$  to denote the set of all distributions over  $X$ . The symbol  $\sqcup$  denotes the disjoint union operation. An *efficient* algorithm is one which runs in probabilistic polynomial time (PPT).

### 2.1 Finding feasible orders for scientific collaboration

We present a model of scientific collaboration, addressing *whether* it is beneficial for parties to collaborate, and if so, *how*: in particular, since the order of research publication matters, we examine the problem of finding a feasible collaboration order.

We propose a model of collaboration between  $n$  research groups which captures the following features. Groups may pool their data, but each group will publish their own results. Moreover, only results that improve on the “state of the art” may be published. That is, a new result must improve on prior publications. However, more credit may be given to earlier publications. Finally, a group will learn not only from pooling their data with other groups, but also from other groups’ publications.

### 2.2 The model

To formalize the intuitions outlined above, we specify a model as follows.

- There is a set  $[n]$  of players.
- Each player  $i \in [n]$  has a dataset  $x_i$  which is sampled in the following way.
  - For each  $i \in [n]$ , there is a set  $X_i$  of possible datasets, which is common knowledge. Let  $X$  denote  $X_1 \times \dots \times X_n$ .
  - There is a distribution  $\mathcal{X} \in \Delta(X)$  over  $X$ , from which the  $x_i$  are sampled:  $(x_1, \dots, x_n) \leftarrow \mathcal{X}$ .
  - The distribution  $\mathcal{X}$  is not known to any of the players, but comes from a commonly known distribution  $\mathcal{D}$ . That is,  $\mathcal{X} \leftarrow \mathcal{D}$ , for some  $\mathcal{D} \in \Delta(\Delta(X))$ .
- There is an output space  $Y$ , and a function  $f : \Delta(X_1 \times \dots \times X_n) \rightarrow Y$  such that  $\hat{y} = f(\mathcal{X})$  is the value which the players wish to learn. That is, the players want to learn some property of the unknown distribution  $\mathcal{X}$  from which their datasets were sampled.  $Y$  and  $f$  are common knowledge.
- $\mathcal{Y}_0$  denotes the distribution of  $\hat{y}$  given  $f$  and  $\mathcal{D}$ .
- There is a *score function*  $s : \Delta(Y) \rightarrow \mathbb{R}_+$ , which varies with  $f$  and  $\mathcal{D}$ . The score function  $s(\cdot)$  is maximized by the distribution  $\hat{\mathcal{Y}}$  which puts probability 1 on the true value  $\hat{y}$ . The score function  $s$  is common knowledge.
  - We require a natural *monotonicity* property of the score function. Namely, let  $\mathcal{Y}$  and  $\mathcal{Z}$  be any distributions, and let  $z$  be a value in the support of  $\mathcal{Z}$ . Then

$$s(\mathcal{Y}) \leq s(\mathcal{Y}|z \leftarrow \mathcal{Z}),$$

where  $z \leftarrow \mathcal{Z}$  denotes the event that  $z$  is sampled from the distribution  $\mathcal{Z}$ .

- *Remark.* Let  $\{\hat{y}|x_1, \dots, x_n\}$  denote the distribution of  $\hat{y}$  given certain datasets  $(x_1, \dots, x_n) \in X$ . A consequence of the monotonicity condition is that given all of the datasets  $x_1, \dots, x_n$  of all players in the model, the best achievable score is  $s(\{\hat{y}|x_1, \dots, x_n\})$ .
- A *collaboration outcome* is given by a permutation  $\pi : [n] \rightarrow [n]$  and a vector of output distributions  $(\mathcal{Z}_1, \dots, \mathcal{Z}_n) \in (\Delta(Y))^n$  such that  $s(\mathcal{Y}_0) < s(\mathcal{Z}_{\pi(1)}) < \dots < s(\mathcal{Z}_{\pi(n)})$ . The intuition behind this condition is that, at time  $t$ , player  $\pi(t)$  will publish  $\mathcal{Z}_{\pi(t)}$ . Since only results that improve on the “state of the art” can be published, we must have that the score  $s(\mathcal{Z}_{\pi(t)})$  increases with the time of publication  $t$ .
- For a collaboration outcome  $\omega = (\pi, \vec{\mathcal{Z}})$ , the player who publishes at time  $t$  obtains a reward

$$R_t(\pi, \vec{\mathcal{Z}}) = \beta^t \cdot (s(\mathcal{Z}_{\pi(t)}) - s(\mathcal{Z}_{\pi(t-1)}))$$

where  $\beta \in (0, 1]$  is a *discount factor* which penalizes later publications.<sup>5</sup>

- For each player  $i$ , we define  $\alpha_i = s(\{\hat{y}|x_i\}) - s(\mathcal{Y}_0) \in \mathbb{R}_+$ , where  $\{\hat{y}|x_i\}$  is the distribution of  $\hat{y}$  given that the  $i^{th}$  dataset is  $x_i$ . This models the “outside payoff” that player  $i$  could get if she does not collaborate and simply publishes on her own.
- Players may learn information not only from their own data, but also from the prior publications of others. A *learning bound vector*  $\{\lambda_{\pi,i}\}_{\pi \in ([n] \rightarrow [n]), i \in [n]}$  characterizes, for any publication order  $\pi$ , the maximum amount that each player  $i$  can learn from prior publications. This notion is defined formally in Section 2.4.
- We define CK to be the collection of all common-knowledge parameters of the model:

$$\text{CK} = (\mathcal{D}, f, s, \beta).$$

## 2.3 Examples

To illustrate the range of settings to which our model can be applied, we now describe several concrete model instantiations.

Recall that our goal is to build mechanisms to enable collaborations by sharing data, in settings where such collaboration would be beneficial to all parties. Intuitively,<sup>6</sup> such settings correspond to situations where the result that can be obtained based on the union of all players’ datasets is “much better” than the results that can be obtained based on the individual datasets: in other words, the “size of the pie” to be split between the players in the case that they collaborate is at least as large as the sum of the “slices” that the players could get by working individually.

### Toy Example I: Secret-sharing

We begin with a “toy example” based on secret-sharing. This first example may seem artificial, but we include it to give a simple initial illustration of the idea that the size of the reward from collaboration can be much greater than what can be achieved without collaborating. We then proceed to three further example settings and explain how our model applies to each.

Consider a stylized secret-sharing model with a secret  $\hat{y}$  drawn uniformly at random from  $\{0, 1\}^n$ . Each player’s data consists of a share  $x_i \in \{0, 1\}^n$  such that  $\hat{y} = x_1 \oplus \dots \oplus x_n$  be the secret the players are trying to reconstruct.

The shares are correlated and drawn from a distribution  $\mathcal{X}$  as follows:

- For each  $i \in [n - 1]$ ,  $x_i$  is uniformly random in  $\{0, 1\}^n$ .
- The last share is chosen such that  $x_n = \hat{y} \oplus x_1 \oplus \dots \oplus x_{n-1}$ .

The players want to learn  $f(\mathcal{X}) = \hat{y}$ . The score from publishing a distribution  $\mathcal{Y}$  is

$$s(\mathcal{Y}) = H(\hat{y}) - H(\hat{y}|\mathcal{Y})$$

where  $H(\hat{y}) = n$  is the entropy of the uniformly random string  $\hat{y}$  and  $H(\hat{y}|\mathcal{Y})$  is the entropy of  $\hat{y}$  given the distribution  $\mathcal{Y}$ .

<sup>5</sup>This is motivated by market scoring rules [Han12], where different experts only receive a reward equal to how much they improve existing predictions.

<sup>6</sup>The intuition we describe here is made rigorous in Claim ??.

Without collaborating, each player  $i$  only knows a uniformly random string  $x_i$ . Thus,  $H(\hat{y}|x_i) = H(\hat{y}) = n$  and  $\alpha_i = H(\hat{y}|x_i) - H(\hat{y}) = 0$  for each player  $i$ . Consider the following collaboration mechanism:

- Each player contributes their share  $x_i$  to the mechanism.
- The mechanism computes  $\hat{y} = x_1 \oplus \dots \oplus x_n$ .
- The mechanism reveals the  $i^{th}$  digit  $\hat{y}_i$  to each player  $i$ .

When participating in this mechanism, the first player will publish a guess  $\mathcal{Y}_1$  which is a distribution over  $\{0, 1\}^n$  where the first bit of  $y \leftarrow \mathcal{Y}_1$  is always  $\hat{y}_1$ . All other players learn  $\hat{y}_1$  from player 1’s publication. Proceeding inductively, the  $i^{th}$  player will publish a guess  $\mathcal{Y}_i$  such that the first  $i$  bits are correct, that is,  $(y_1, \dots, y_i) = (\hat{y}_1, \dots, \hat{y}_i)$  for any  $y \leftarrow \mathcal{Y}_i$ . Note that since  $\alpha_i = 0$  for each player  $i$ , and  $H(\hat{y}|\mathcal{Y}_i) - H(\hat{y}|\mathcal{Y}_{i-1}) = 1 > \alpha_i$ , this mechanism incentivizes players to collaborate.

### Toy Example II: Network flow

Let  $G = (V, E)$  be a graph where  $V$  is the set of vertices and  $E$  is the set of edges. Let  $\tilde{s}, \tilde{t} \in V$  be vertices which are connected by some number of disjoint paths. Consider a model where  $V$ ,  $\tilde{s}$ , and  $\tilde{t}$  are common knowledge, and each player’s data consists of a disjoint subset of edges in  $x_i \subseteq E$ . More precisely,  $(x_1, \dots, x_n) \leftarrow \mathcal{X}(E)$  where  $\mathcal{X}$  samples a partition of  $E$ .

The players want to learn the set of paths from  $\tilde{s}$  to  $\tilde{t}$ . That is,  $f(\mathcal{X}(E))$  is the set of paths in  $E$  that start at  $\tilde{s}$  and end at  $\tilde{t}$ . The score from publishing a distribution  $\mathcal{Z}$  over edges is

$$s(\mathcal{Z}) = |\{p : p \text{ is a path in } E \text{ from } \tilde{s} \text{ to } \tilde{t}, \text{ and } \Pr_{z \leftarrow \mathcal{Z}} [p \subseteq z] = 1\}|.$$

In other words, the player’s score is given by how many paths from  $\tilde{s}$  to  $\tilde{t}$  she knows with certainty to exist in  $E$ . In some cases, it may be that no player knows any path from  $\tilde{s}$  to  $\tilde{t}$  based only on her own data, as illustrated by the simple example in Figure 1.

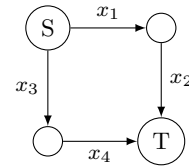


Figure 1: A four-player collaboration example

Consider the following collaboration mechanism:

- Each player contributes their set of edges  $x_i$  to the mechanism.
- The mechanism computes  $E = x_1 \cup \dots \cup x_n$ , and then computes the set  $P = \{p_1, \dots, p_k\}$  of paths in  $E$  that start at  $\tilde{s}$  and end at  $\tilde{t}$ .

- The mechanism reveals the  $i^{th}$  path  $p_i$  to player  $i$ . If  $k < n$ , then the last  $k - n$  players will get no output. If  $k > n$ , the “extra” paths are allocated arbitrarily to players.<sup>7</sup>

When participating in this mechanism, the first player will publish a guess  $\mathcal{Z}_1$  which (always) samples the set  $\{p_1\}$ . All other players learn  $p_1$  from player 1’s publication. Then, the  $i^{th}$  player will publish a guess  $\mathcal{Z}_i$  that samples the set  $\{p_1, \dots, p_i\}$ . As long as  $s(\mathcal{Z}_i) - s(\mathcal{Z}_{i-1}) \geq \alpha_i$  for all  $i \in [n]$  (note that this is the case in Figure 1), this mechanism incentivizes players to collaborate.

### Example III: Correlating gene loci with disease

This example is inspired by successful GWAS studies to identify gene loci associated with schizophrenia. Consider a model where each player holds a set of patients’ medical (and in particular, genetic) data  $x_i$  which comes from some unknown patient distribution  $\mathcal{X}$ . The players wish to learn the set  $f(\mathcal{X})$  of gene loci that are correlated with the occurrence of schizophrenia in patients.

Let  $\Gamma$  be the set of all gene loci. For  $\gamma \in \Gamma$ , define  $\mathbb{I}_\gamma$  to be 1 if  $\gamma \in f(\mathcal{X})$  and 0 otherwise. The score from publishing a distribution  $\mathcal{Z}$  over  $\mathcal{P}(\Gamma)$  (i.e. over subsets of gene loci) could be:<sup>8</sup>

$$s(\mathcal{Z}) = \sum_{\gamma \in f(\mathcal{X})} \Pr[\gamma \in \mathcal{Z}] - \sum_{\gamma \notin f(\mathcal{X})} \Pr[\gamma \in \mathcal{Z}].$$

This score function rewards players for assigning high probabilities to gene loci  $\gamma$  which are actually correlated with schizophrenia, and penalizes them for assigning high probabilities to those which are not. As in our previous examples, it turns out that in this setting, the reward that can be obtained based on pooling all the players’ data is much greater than the sum of the rewards that could be obtained individually, as illustrated in Figure 2.

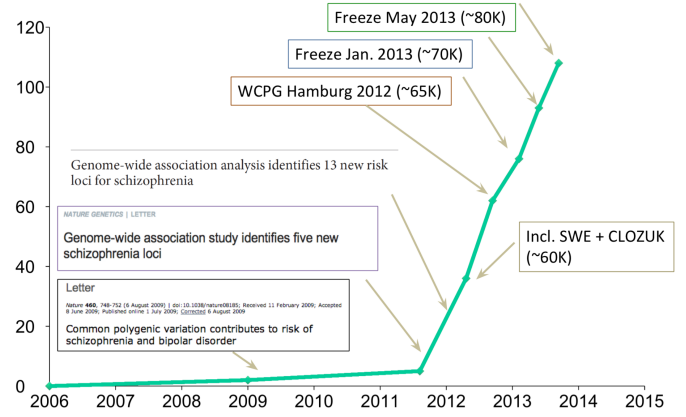
Consider the following collaboration mechanism:<sup>11</sup>

<sup>7</sup>It may be beneficial to allocate the “extra” paths strategically in order to reward players more fairly, or in order to make collaboration possible when the outside option values  $\alpha_i$  are nonzero. However, in this example, we allocate them arbitrarily for simplicity.

<sup>8</sup>In practice, a more realistic scenario might be to model the *extent* to which particular gene loci are found to be correlated with the occurrence of schizophrenia, rather than classifying into binary categories “correlated” and “not correlated”. This case could be modeled, for example, by letting  $f(\mathcal{X})$  be a vector  $((\gamma_1, p_1), \dots, (\gamma_N, p_N))$  where  $\Gamma = \{\gamma_1, \dots, \gamma_N\}$  is the set of gene loci, and for each  $j \in [N]$ ,  $p_j$  is the correlation coefficient between  $\gamma_j$  and occurrence of schizophrenia. While Example IV presents the simpler “binary” model for ease of exposition, we remark that with appropriate modifications to the score function and mechanism, our model can accommodate the more complex case of estimating correlations, too.

<sup>10</sup>Graph taken from a class given by Manolis Kellis; GWAS studies performed by the Schizophrenia Genomics Consortium.

<sup>11</sup>This is just one example of a reasonable mechanism for this model; we do not mean to claim that it is a canonical or optimal one. There are many variants which could make sense: for example, a simple modification would be to change the threshold 0.5 used in the second step.



**Figure 2: Graph of GWAS study success: the y-axis is the number of gene loci identified to be correlated with schizophrenia, and the x-axis is time (which corresponds to *amount of data*, since the reason for the improved findings was accumulation of data over time).<sup>10</sup>**

- Each player contributes their patient data  $x_i$  to the mechanism.
- The mechanism computes  $\mathcal{Y}^* = \{f(\mathcal{X})|x_1, \dots, x_n\}$ , i.e. the distribution of  $f(\mathcal{X})$  given the input data of all the players. Let  $\Gamma^* = \{\gamma \in \Gamma : \Pr_{y \leftarrow \mathcal{Y}^*}[\gamma \in y] > 0.5\}$ , that is, the set of gene loci that are more likely than not to be in  $f(\mathcal{X})$ , according to  $\mathcal{Y}^*$ .
- The mechanism reveals to player  $i$  the  $i^{th}$  gene locus  $\gamma_i$  in  $\Gamma^*$ . If  $|\Gamma^*| < n$ , then the last  $k - n$  players will get no output. If  $|\Gamma^*| > n$ , the “extra” gene loci are allocated arbitrarily.<sup>12</sup>

When participating in this mechanism, the first player will publish a guess  $\mathcal{Z}_1$  which (always) samples the set  $\{\gamma_1\}$ . All other players learn  $\gamma_1$  from player 1’s publication. Then, the  $i^{th}$  player will publish a guess  $\mathcal{Z}_i$  that samples the set  $\{\gamma_1, \dots, \gamma_i\}$ . Provided that  $s(\mathcal{Z}_i) - s(\mathcal{Z}_{i-1}) \geq \alpha_i$  for all  $i \in [n]$  (note that Figure 2 depicts exactly such a scenario), this mechanism incentivizes players to collaborate.

### Example IV: Statistical estimation

Our last example is one where – in contrast to the examples so far – there are *decreasing* marginal returns from adding new information, and thus collaboration will not be feasible.

We consider a simple Bayesian model where the distribution  $\mathcal{X}$  is itself drawn from a “distribution over distributions”  $\mathcal{D}$ . More concretely, each player  $i$  receives a vector of  $k_i$  samples  $(x_{i,1}, \dots, x_{i,k_i})$  drawn independently from a normal distribution  $N(\mu, \sigma^2)$  with unknown mean  $\mu$  and known variance  $\sigma^2$ . The mean  $\mu$  is itself drawn from a commonly known prior distribution  $\mathcal{D} = N(m, 1)$  with known mean  $m$  and variance 1. In this case, the ground set  $X_i$  is  $\mathbb{R}^{k_i}$ . The distribution  $\mathcal{X}(\mu, \sigma)$  is a product distribution

<sup>12</sup>As remarked in Footnote 7, it can be beneficial to allocate the “extra” gene loci in a way which is not arbitrary, but instead optimized for making collaboration possible. In this example, for simplicity, we allocate them arbitrarily.

over  $\mathbb{R}^{\sum_{i=1}^n k_i}$ , where each component of  $(x_{1,1}, \dots, x_{n,k_n})$  is drawn independently from  $N(\mu, \sigma)$ . The players want to learn  $f(\mathcal{X}(\mu, \sigma)) = \mu$ .

An estimator for  $\mu$  is a random variable  $\hat{\mu}$ . The score of such a guess  $\hat{\mu}$  is  $s(\hat{\mu}) = -\mathbb{E}[(\hat{\mu} - \mu)^2]$ . It is well known that if we have a vector  $(x_{i,1}, \dots, x_{i,k_i})$  of random samples drawn from  $N(\mu, \sigma)$ , the estimator that minimizes the expected squared error to  $\mu$  is  $\hat{\mu}_i = \frac{1}{k_i} \sum_{j=1}^{k_i} x_{i,j}$ . Note that this is a normal random variable since each  $x_{i,j}$  is sampled from normal random variable. The expectation of  $\hat{\mu}_i$  is  $\frac{1}{k_i} \cdot k_i \cdot \mu = \mu$  and the variance of  $\hat{\mu}_i$  is  $\frac{1}{k_i^2} \cdot k_i \cdot \sigma^2 = \frac{1}{k_i} \cdot \sigma^2$ . Thus,  $s(\hat{\mu}_i) = \frac{1}{k_i} \cdot \sigma^2$ . If a player published by herself and did not collaborate, her reward would be the difference  $\alpha_i = \sigma^2 - \frac{1}{k_i} \cdot \sigma^2$  between the priorly known variance  $\sigma^2$  and the variance  $\frac{1}{k_i} \cdot \sigma^2$  of player  $i$ 's estimate.

If the players collaborate, they can obtain the estimator  $\hat{\mu}^* = \frac{1}{\sum_{i=1}^n k_i} \sum_{i=1}^n \sum_{j=1}^{k_i} x_{i,j}$  which has variance  $s(\hat{\mu}^*) = \frac{1}{\sum_{i=1}^n k_i} \sigma^2$ . The reward for  $\hat{\mu}^*$  is the reduction in variance  $\sigma^2 - s(\hat{\mu}^*) = \sigma^2 \cdot (1 - \frac{1}{\sum_{i=1}^n k_i})$ . Note that in this case, the reward from an estimator only depends on the number of data points  $N$  used to construct this estimator (in the above notation,  $N = \sum_{i=1}^n k_i$ ). Furthermore, the reward  $R(N) = \sigma^2(1 - \frac{1}{N})$  that one could obtain with  $N$  data points is concave in  $N$ . Intuitively, if one only has  $N = 2$  data points, and gets 10 new ones, those 10 new data points are very valuable. However, if one already has  $N = 2000000$  data points and gets 10 new ones, those 10 new data points do not increase the score very much.

This setting is in contrast to our Example III, where the score seemed to increase in a convex way with the number of data points. Indeed, in this Bayesian example, we will always have that

$$R(\sum_{i=1}^n k_i) = \sigma^2(1 - \frac{1}{\sum_{i=1}^n k_i}) \leq \sigma^2 \sum_{i=1}^n (1 - \frac{1}{k_i}) = \sum_{i=1}^n R(k_i).$$

We will elaborate more formally in Section 2.5 why the above inequality is bad for collaboration, but we can already see intuitively from it that the left-hand side is the "size of the pie" if all players were to collaborate, and the right-hand side is the sum of the rewards that each player could receive on her own. Because of the above inequality, there is no way to "slice the pie" so that every player has a bigger reward than the  $\alpha_i$  they can get without collaborating, and thus collaboration is impossible.

This raises the interesting question of *when* the value of information is convex with the amount of information available. From the previous three examples, we have seen that this is often the case. However, from this example, we see that even in a simple Bayesian setting, the marginal value of extra information will be decreasing. Thus, this suggests that collaboration will be more feasible when learning is done in a non-Bayesian way.

## 2.4 Data-sharing mechanisms

We now return to the general formulation of our collaboration model, and we seek to design a general data-sharing mechanism that takes as input the data of all the parties, computes an output distribution  $\mathcal{Y}_i \in \Delta(Y)$  for each  $i \in [n]$ , and outputs  $\mathcal{Y}_i$  to each player  $i$ . The mechanism will output the  $\mathcal{Y}_i$  values to players sequentially, in a particular order.

Upon receiving  $\mathcal{Y}_i$ , player  $i$  produces a public output (i.e a publication in the research collaboration example) which we denote by  $\mathcal{Z}_i \in Y$ .

We note that the public output of player  $i$  will not necessarily be the same as what was delivered by the data-sharing mechanism. Since player  $i$  wants to maximize her reward, she will publish a result  $\mathcal{Z}_i$  that will maximize her reward, conditional on the information she has at the time of publication. This information includes, in addition to the output  $\mathcal{Y}_i$  which she receives from the mechanism (and her knowledge of how the mechanism works<sup>13</sup>), also her own dataset  $x_i \in X_i$ , and all the outputs  $\mathcal{Z}_j$  of other players that published before her.

Recall that a *collaboration outcome*  $(\pi, \vec{\mathcal{Z}})$  is given by a permutation  $\pi : [n] \rightarrow [n]$  and a vector of output distributions  $\vec{\mathcal{Z}} = (\mathcal{Z}_1, \dots, \mathcal{Z}_n) \in (\Delta(Y))^n$  such that  $s(\mathcal{Y}_0) < s(\mathcal{Z}_{\pi(1)}) < \dots < s(\mathcal{Z}_{\pi(n)})$ . We now define a *proposed collaboration outcome*  $(\pi, \vec{\mathcal{Y}})$  as a permutation  $\pi : [n] \rightarrow [n]$  together with a vector of *proposed* outputs  $\vec{\mathcal{Y}} = (\mathcal{Y}_1, \dots, \mathcal{Y}_n) \in (\Delta(Y))^n$  generated by a data-sharing mechanism, satisfying  $s(\mathcal{Y}_0) < s(\mathcal{Y}_{\pi(1)}) < \dots < s(\mathcal{Y}_{\pi(n)})$ .

Recall also that we need to bound how much player  $i$  can learn from previous publications (and from her own dataset). We formally capture this with the notion of *learning bound vectors*  $\lambda_{\pi,i}$ , which give an upper bound on the amount that player  $i$  learns from all previous publications when the order of publication is determined by permutation  $\pi$ .

DEFINITION 1. A learning bound vector

$$\vec{\lambda} = (\lambda_{\pi,i})_{\pi \in ([n] \rightarrow [n]), i \in [n]}$$

is a non-negative vector such that, if  $(\pi, \vec{\mathcal{Y}})$  is a collaboration outcome proposed by a data-sharing mechanism, and  $\mathcal{Z}_i$  is the best (i.e. highest-scoring) distribution that player  $i$  can compute at the time  $\pi^{-1}(i)$  of her publication, then  $s(\mathcal{Z}_i) \leq s(\mathcal{Y}_i) + \lambda_{\pi,i}$ . Let  $\Lambda = \mathbb{R}_+^{n! \times n}$  denote the set of all learning bound vectors.

DEFINITION 2. For a learning bound vector  $\vec{\lambda}$ , the set of inferred output distributions derived from a proposed collaboration outcome  $(\pi, \vec{\mathcal{Y}})$  is given by the following expression:

$$\mathcal{I}_{\vec{\lambda}}(\pi, \vec{\mathcal{Y}}) = \{(\mathcal{Z}_1, \dots, \mathcal{Z}_n) : \forall t \in [n], s(\mathcal{Y}_{\pi(t)}) \leq s(\mathcal{Z}_{\pi(t)}) \leq s(\mathcal{Y}_{\pi(t)}) + \lambda_{\pi,\pi(t)}\}.$$

The intuition behind the above definition is that the amount of information that player  $\pi(t)$  (namely, the player who publishes at time  $t$ ) can learn from prior outputs is measured by how much her score increases based on these prior outputs. This increase in score is bounded by  $\lambda_{\pi,\pi(t)}$ . Thus, her eventual output will be some  $\mathcal{Z}_{\pi(t)}$  with score between  $s(\mathcal{Y}_{\pi(t)})$  and  $s(\mathcal{Y}_{\pi(t)}) + \lambda_{\pi,\pi(t)}$ .

REMARK 1. In certain cases,  $\lambda_{\pi,\pi(t)}$  measures exactly the amount of information that player  $\pi(t)$  can learn from her data. However, in our definition  $\lambda_{\pi,\pi(t)}$  is an upper bound, and we emphasize that it may be a loose upper bound on the amount of information  $\pi(t)$  can learn. Our emphasis on this point comes from the following two reasons.

- In general, the vector  $\vec{\lambda} \in \mathbb{R}^{n! \times n}$  has very high dimension, and finding such a vector is infeasible. We may

<sup>13</sup>The mechanism description is common knowledge.



want to approximate this vector via a low-dimensional encoding (as we will do below, where we encode learning bounds using  $n$ -dimensional vectors). Since this low-dimensional encoding will lose information, we will not be able to represent  $\lambda_{\pi, \pi(t)}$  exactly, but may get a reasonable upper bound on its value.

- For some other settings, we may not be able to derive a precise expression for  $\lambda_{\pi, \pi(t)}$  in terms of expectations, but we may still be able to derive an upper bound on the amount of information that player  $\pi(t)$  learns. An example of such a setting is detailed in Section ??.

Now that we have established a formal definition of learning bound vectors, we proceed to formally define a data-sharing mechanism.

DEFINITION 3. For model parameters CK, a data sharing mechanism is a function

$$M : X \times \Lambda \rightarrow ([n] \rightarrow [n]) \times (\Delta(Y))^n$$

which takes as inputs a vector  $\vec{x} = (x_1, \dots, x_n)$  of datasets and  $\vec{\lambda} = (\lambda_{\pi, i})_{\pi \in ([n] \rightarrow [n]), i \in [n]}$  a learning bound vector, and outputs an ordering  $\pi$  of the players and an output vector  $(\mathcal{Y}_1, \dots, \mathcal{Y}_n) \in (\Delta(Y))^n$ .

REMARK 2. In the definition, for the sake of generality, we assume that the  $\vec{\lambda}$  values are given as input to the mechanism. We remark that in certain settings, these values can be computed directly from the inputs  $x_i$  of the parties, as discussed in the examples of Section 1.1.2. In this case, one may think of the mechanism  $M : X \rightarrow ([n] \rightarrow [n]) \times (\Delta(Y))^n$  as having input domain  $X$  only.

## 2.5 Collaborative equilibria

In our model, each research group  $\pi(t)$  will collaborate only if the credit they obtain from doing so is greater than the “outside option” reward  $\alpha_{\pi(t)}$ . We want to design a mechanism that guarantees collaboration whenever possible. Accordingly, we define the following equilibrium concept.

DEFINITION 4. Let CK be the model parameters. Let  $(\vec{x}, \vec{\lambda}) \in X \times \Lambda$  and let  $(\pi, (\mathcal{Y}_1, \dots, \mathcal{Y}_n)) \in ([n] \rightarrow [n]) \times (\Delta(Y))^n$ . We say that  $(\pi, (\mathcal{Y}_1, \dots, \mathcal{Y}_n))$  is a collaborative equilibrium with respect to  $(\vec{x}, \vec{\lambda})$  if for all inferred output distributions  $\vec{\mathcal{Z}} = (\mathcal{Z}_1, \dots, \mathcal{Z}_n) \in \mathcal{I}(\pi, (\mathcal{Y}_1, \dots, \mathcal{Y}_n))$  and all  $t \in [n]$ , it holds that  $R_t(\pi, \vec{\mathcal{Z}}) \geq \alpha_{\pi(t)}$ .

Our goal is to find data-sharing mechanisms for which collaboration is an equilibrium. Intuitively, since we are searching for a feasible permutation over a very high-dimensional space ( $n!$ -dimensional, to be precise), the problem will be NP-complete (this is proven in Theorem 3). However, there is a very natural condition on the learning vectors for which we can reduce the dimension of the search space and efficiently find a collaborative equilibrium. The feasible case corresponds to the case where, for any player  $j$ , there is a bound on the amount of information that player  $j$  could teach any other players. We denote this bound by  $\mu_j$ . Analogously, we could define  $\mu_j$  to be a bound on the amount that player  $j$  can learn from any other player. In this work, we describe only the first case, when  $\mu_j$  represents a bound on how much information player  $j$  can teach other players. The other case is analogous.

We define a learning bound vector to be  $n$ -dimensional if it satisfies the following property.

DEFINITION 5. A learning vector  $\vec{\lambda} \in \Lambda$  is  $n$ -dimensional if there is a non-negative vector  $(\mu_1, \dots, \mu_n)$  such that  $\lambda_{\pi, \pi(t)} = \sum_{\tau=1}^{t-1} \mu_{\pi(\tau)}$ . Let  $\Lambda_1 \subset \Lambda$  denote the set of all  $n$ -dimensional learning vectors.

When  $\vec{\lambda}$  is an  $n$ -dimensional learning vector, the total amount that player  $\pi(t)$  learns from all prior outputs is  $\sum_{\tau=1}^{t-1} \mu_{\pi(\tau)}$ . In this case, we can give necessary and sufficient conditions for an equilibrium to exist (detailed in Theorem 1 below), provided that the following Output Divisibility Condition is satisfied.

### Output Divisibility Condition..

Given the model parameters CK and any real  $0 < \delta \leq 1$ ,<sup>14</sup> there exists a distribution  $\mathcal{Y} \in \Delta(Y)$  such that  $s(\mathcal{Y}) = \delta$ .

REMARK 3. The Output Divisibility Condition holds for a wide variety of natural score functions. In general, score functions which reward “how close” a distribution is to the true value  $\hat{y} = f(\mathcal{X})$  will decrease (continuously) with the addition of random noise to a distribution. Provided that this holds, the Output Divisibility Condition can be satisfied by taking the optimal distribution  $\{\hat{y}|\mathcal{X}\}$  and perturbing it with random noise: the exact amount of noise to be added depends on the desired value of  $\delta$ . To give a concrete example: in Example III (Gene loci), the perturbed distribution could simply add noise to the probabilities that each gene locus is sampled. Here, “adding noise” can mean simply adding some  $\eta \leftarrow N(0, \sigma^2)$  to the relevant parameters, where the magnitude of  $\sigma$  depends on the precise formulation of the score function and the desired value of  $\delta$ .

THEOREM 1. Suppose that the Output Divisibility Condition holds. Let  $\vec{x}$  be a vector of inputs and  $\vec{\lambda}$  be an  $n$ -dimensional learning bound vector. Let  $\lambda_{\pi, \pi(t)} = \sum_{\tau=1}^{t-1} \mu_{\pi(\tau)}$ . Then for  $(\pi, \vec{\mathcal{Y}})$  to be a collaborative equilibrium, it is necessary and sufficient that

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n (n-t) \mu_{\pi(t)} \leq s(\mathcal{Y}_{\pi(n)}) - s(\mathcal{Y}_0).$$

PROOF. **Necessity.** Let  $(\pi, \vec{\mathcal{Y}})$  be a proposed collaborative equilibrium, and let  $\vec{\mathcal{Z}} \in \mathcal{I}(\pi, \vec{\mathcal{Y}})$  be a possible vector of inferred outputs. For every  $t$ , we must have that:

$$\beta^t \cdot (s(\mathcal{Z}_{\pi(t)}) - s(\mathcal{Z}_{\pi(t-1)})) \geq \alpha_{\pi(t)}.$$

This is equivalent to:

$$s(\mathcal{Z}_{\pi(t)}) - s(\mathcal{Z}_{\pi(t-1)}) \geq \frac{\alpha_{\pi(t)}}{\beta^t}.$$

The worst case for player  $\pi(t)$  is when player  $\pi(t-1)$  learns as much as possible from prior publications and player  $\pi(t)$  learns as little as possible. That is, when

$$s(\mathcal{Z}_{\pi(t-1)}) = s(\mathcal{Y}_{\pi(t-1)}) + \mu_{\pi(1)} + \dots + \mu_{\pi(t-2)} \\ \text{and } s(\mathcal{Z}_{\pi(t)}) = s(\mathcal{Y}_{\pi(t)}).$$

<sup>14</sup>Recall (from the model description) that  $s(\{\hat{y}|\mathcal{X}\}) = \max_{\mathcal{Y} \in \Delta(Y)} (s(\mathcal{Y}))$ . Without loss of generality, we assume in our analysis that the score function is normalized so that its maximum value  $s(\{\hat{y}|\mathcal{X}\}) = 1$ .

In this case, the equilibrium condition becomes:

$$s(\mathcal{Y}_{\pi(t)}) - s(\mathcal{Y}_{\pi(t-1)}) - \sum_{\tau=1}^{t-2} \mu_{\pi(\tau)} \geq \frac{\alpha_{\pi(t)}}{\beta^t}.$$

Rearranging slightly, we obtain:  $s(\mathcal{Y}_{\pi(t-1)}) - s(\mathcal{Y}_{\pi(t)}) \leq -\frac{\alpha_{\pi(t)}}{\beta^t} - \sum_{\tau=1}^{t-2} \mu_{\pi(\tau)}$ . Let us abuse notation slightly and define  $\pi(0) = 0$ . Then, summing over all  $t$  yields

$$s(\mathcal{Y}_{\pi(0)}) - s(\mathcal{Y}_{\pi(n)}) \leq -\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} - \sum_{t=1}^n (n-t) \mu_{\pi(t)}.$$

Flipping the signs in the inequality, the existence of a collaborative equilibrium implies:

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n (n-t) \mu_{\pi(t)} \leq s(\mathcal{Y}_{\pi(n)}) - s(\mathcal{Y}_0).$$

**Sufficiency.** To prove that the condition is sufficient: given  $\mathcal{Y}_{\pi(n)}$  satisfying the inequality in the theorem statement, we need to construct  $\vec{\mathcal{Y}} = (\mathcal{Y}_1, \dots, \mathcal{Y}_n)$  such that  $(\pi, \vec{\mathcal{Y}})$  is a collaborative equilibrium. We construct  $\vec{\mathcal{Y}}$  inductively as follows: let  $\delta_{\pi(n)} = s(\mathcal{Y}_{\pi(n)})$ , and for any  $t$  such that  $2 \leq t \leq n$ , let  $\delta_{\pi(t-1)} = \delta_{\pi(t)} - \frac{\alpha_{\pi(t)}}{\beta^t} - \sum_{\tau=1}^{t-2} \mu_{\pi(\tau)}$ . Now that we have defined  $\{\delta_{\pi(t)}\}_{t=1}^n$  in this way, it follows that if we set  $\mathcal{Y}_{\pi(t)}$  such that  $s(\mathcal{Y}_{\pi(t)}) = \delta_{\pi(t)}$ , then for all  $t \geq 2$  we have

$$s(\mathcal{Y}_{\pi(t)}) - s(\mathcal{Y}_{\pi(t-1)}) = \delta_{\pi(t)} - \delta_{\pi(t-1)} = \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{\tau=1}^{t-2} \mu_{\pi(\tau)}.$$

Note that it is possible to set  $\mathcal{Y}_{\pi(t)}$  in the required way, by the Output Divisibility Condition. Rearranging the above equation, it follows that:

$$\beta^t \cdot (s(\mathcal{Y}_{\pi(t)}) - s(\mathcal{Y}_{\pi(t-1)}) - \sum_{\tau=1}^{t-2} \mu_{\pi(\tau)}) = \alpha_{\pi(t)}.$$

Since for any inferred outcome  $\mathcal{Z}_{\pi(t)}$  we have (by the definition of the learning bound vector) that

$$s(\mathcal{Y}_{\pi(t)}) \leq s(\mathcal{Z}_{\pi(t)}) \leq s(\mathcal{Y}_{\pi(t)}) + \lambda_{\pi, \pi(t-1)} = s(\mathcal{Y}_{\pi(t)}) + \sum_{\tau=1}^{t-2} \mu_{\pi(\tau)},$$

we conclude that for all  $t \geq 2$ ,

$$\beta^t \cdot (s(\mathcal{Z}_{\pi(t)}) - s(\mathcal{Z}_{\pi(t-1)})) \geq \alpha_{\pi(t)}.$$

Finally, we need to check that player  $\pi(1)$  is incentivized to collaborate. Note that player  $\pi(1)$  publishes first, so she cannot learn anything from previous publications. She will be incentivized to publish if

$$\beta \cdot (\delta_{\pi(1)} - s(\mathcal{Y}_0)) \geq \alpha_{\pi(1)}.$$

This condition is equivalent to

$$\delta_{\pi(1)} - s(\mathcal{Y}_0) \geq \frac{\alpha_{\pi(1)}}{\beta}.$$

Replacing  $\delta_{\pi(t-1)} = \delta_{\pi(t)} - \frac{\alpha_{\pi(t)}}{\beta^t} - \sum_{\tau=1}^{t-2} \mu_{\pi(\tau)}$  iteratively, we get that player  $\pi(1)$  is incentivized to collaborate if and only if

$$\delta_{\pi(n)} - s(\mathcal{Y}_0) \geq \sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n (n-t) \mu_{\pi(t)}$$

which is guaranteed by assumption.  $\square$

Recall from the definition of the score function that the best score that can be attained given datasets  $x_1, \dots, x_n$  is equal to  $s(\{\hat{y}|x_1, \dots, x_n\})$ . Based on Theorem 1, we can now characterize the datasets and learning bound vectors for which a collaborative equilibrium is possible.

**DEFINITION 6.** Let  $\mathbf{CK}$  be the model parameters and let  $(\vec{x}, \vec{\lambda}) \in X \times \Lambda$ . We say that  $(\vec{x}, \vec{\lambda})$  supports a collaborative equilibrium if it holds that

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n (n-t) \mu_{\pi(t)} \leq s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0).$$

### 2.5.1 How do the model parameters affect feasibility of collaborative equilibria?

Consider for a moment the simple case where  $\beta = 1$  and  $\vec{\lambda} = \vec{0}$ , that is, there is no discount factor and players do not learn from others' publications. We can show that in this case, if the score function satisfies the following Property 1, then it holds that for all  $\vec{x} \in X$ ,  $(\vec{x}, \vec{\lambda})$  supports a collaborative equilibrium. That is, in this simple case, the condition for  $(\vec{x}, \vec{\lambda})$  to support an equilibrium reduces to the superadditivity of the auxiliary score function  $\bar{s}$  given in Property 1.

**DEFINITION 7.** Let  $S$  be a set. A function  $f : S \rightarrow \mathbb{R}$  is superadditive if for all disjoint  $S_1, S_2 \subseteq S$ , it holds that  $f(S_1) + f(S_2) \leq f(S_1 \cup S_2)$ .

**PROPERTY 1 (SUPERADDITIVE DIFFERENCES).** Let  $\mathbf{CK}$  be the model parameters. We define an auxiliary score function  $\bar{s} : X_1 \sqcup \dots \sqcup X_n \rightarrow \mathbb{R}_+$  which maps a set of datasets to a real-valued score, as follows:

$$\bar{s}(\{(i_1, x_{i_1}), \dots, (i_k, x_{i_k})\}) = s(\{\hat{y}|x_{i_1}, \dots, x_{i_k}\}) - s(\mathcal{Y}_0),$$

where  $\{\hat{y}|x_{i_1}, \dots, x_{i_k}\}$  denotes the distribution of  $\hat{y}$  given that the datasets  $x_{i_1}, \dots, x_{i_k}$  were sampled<sup>15</sup> from  $\mathcal{X}$ . The score function  $s$  satisfies the Superadditive Differences Property if  $\bar{s}$  is superadditive.

We observe that this precisely captures the intuition initially described in Section 2.3, that our model is designed to promote collaboration in situations where the reward that can be obtained from pooling all players' data is more than the sum of the individual rewards that players can get.

**CLAIM 1.** Let  $\mathbf{CK}$  be model parameters such that  $\beta = 1$ , let  $\vec{x} \in X$  be arbitrary, and let  $\vec{\lambda} = \vec{0} \in \Lambda$ . If  $\bar{s}$  is a superadditive function on the input data, then  $(\vec{x}, \vec{\lambda})$  supports a collaborative equilibrium.

**PROOF.** Recall the inequality from Definition 6:

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n (n-t) \mu_{\pi(t)} \leq s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0).$$

Since  $\beta = 1$  and  $\mu_{\pi(t)} = 0$ , the left-hand side is simply  $\sum_{t=1}^n \alpha_{\pi(t)}$ . Using the definitions of  $\alpha_{\pi(t)}$  and  $\bar{s}$ , and the

<sup>15</sup>More precisely:  $\{\hat{y}|x_{i_1}, \dots, x_{i_k}\}$  is the distribution of  $\hat{y}$  given that each  $x_{i_j}$  was sampled in the  $i_j^{th}$  position. (Recall that the distribution  $\mathcal{X}$  is over tuples of datasets  $(x_1, \dots, x_n)$ .)

fact that  $\pi$  is a permutation, this can be rewritten as:

$$\sum_{t=1}^n \alpha_{\pi(t)} = \sum_{t \in [n]} (s(\{\hat{y}|x_{\pi(t)}\}) - s(\mathcal{Y}_0)) = \sum_{i \in [n]} \bar{s}(\{(i, x_i)\}).$$

Substituting back into the inequality, we obtain:

$$\sum_{i \in [n]} \bar{s}(\{(i, x_i)\}) \leq s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0).$$

The right-hand side of the inequality is, by definition, equal to  $\bar{s}(\{(1, x_1), \dots, (n, x_n)\})$ . Thus, the superadditivity of  $\bar{s}$  implies that the inequality holds, and it follows that  $(\vec{x}, \vec{\lambda})$  supports a collaborative equilibrium.  $\square$

Finally, we remark that either decreasing the discount factor  $\beta$  or increasing the learning bound vector  $\vec{\lambda}$  will make it harder to support a collaborative equilibrium (i.e. a lower value of  $\beta$  means there will be fewer  $(\vec{x}, \vec{\lambda})$  which support an equilibrium), since these cause the left-hand side of the inequality to increase. So, while superadditivity is a sufficient condition in the simplest case, we observe that determining which  $(\vec{x}, \vec{\lambda})$  support a collaborative equilibrium is a more complex problem when the model parameters are varied.

## 2.6 The polynomial-time mechanism

We show a polynomial-time mechanism that computes a collaborative equilibrium in the case that learning bounds are given by a  $n$ -dimensional vector, provided that the following *Efficient Output Divisibility Condition* is satisfied. The *Efficient Output Divisibility Condition* is a natural extension of the *Output Divisibility Condition*, which requires not only existence but also efficient computability of distributions with arbitrary score, while taking into account that the best possible score for given input datasets  $x_1, \dots, x_n$  is equal to  $s(\{\hat{y}|x_1, \dots, x_n\})$ .

### *Efficient Output Divisibility Condition..*

Given model parameters CK, datasets  $x_1, \dots, x_n \in X$ , and any real  $0 < \delta < s(\{\hat{y}|x_1, \dots, x_n\})$ , it is possible to efficiently compute a distribution  $\mathcal{Y} \in \Delta(Y)$  such that  $s(\mathcal{Y}) = \delta$ .

**REMARK 4.** *The Efficient Output Divisibility Condition holds for a wide variety of score functions, too: in particular, it holds for the natural class of score functions described in Section 2.5 as satisfying the Output Divisibility Condition. Suppose that the score function is continuous and decreases with the addition of random noise to a distribution. Then the Efficient Output Divisibility Condition can be satisfied by taking the “best computable” distribution  $\{\hat{y}|x_1, \dots, x_n\}$  and perturbing it with random noise: the exact amount of noise to add will depend on the desired value of  $\delta$ . As discussed in the context of the Output Divisibility Condition, this can also be easily explained in terms of the concrete Examples III and IV.*

**THEOREM 2.** *Suppose the Efficient Output Divisibility Condition holds. Then there is a polynomial-time mechanism SHARE-DATA :  $X \times \Lambda_1$  that, given inputs  $(\vec{x}, \vec{\mu})$  where  $\vec{\mu} = (\mu_1, \dots, \mu_n)$  represents a  $n$ -dimensional learning vector, outputs a collaborative equilibrium  $(\pi, \vec{\mathcal{Y}})$  whenever an equilibrium is supported by the inputs  $(\vec{x}, \vec{\mu})$  (as defined in Definition 6), and outputs NONE otherwise.*

---

### Algorithm 1 SHARE-DATA( $(x_1, \dots, x_n), (\mu_1, \dots, \mu_n)$ )

---

1. Let  $\mathcal{Y}^* = \{\hat{y}|x_1, \dots, x_n\}$  and  $\delta^* = s(\mathcal{Y}_0)$ .
  2. Construct a complete weighted bipartite graph  $G = (L, R, E)$  where  $L = [n], R = [n], E = L \times R$ . For each edge  $(i, t)$ , assign a weight  $w(i, t) = \frac{\alpha_i}{\beta^t} + (n - t)\mu_i$ .
  3. Let  $M$  be the minimum-weight perfect matching on  $G$ . For each node  $t \in R$ , let  $\pi(t) \in L$  be the node that it is matched with. If the weight of  $M$  is larger than  $\delta^*$ , output NONE. Else, define  $\delta_{\pi(n)} = \delta^*$ ,  $\mathcal{Y}_{\pi(n)} = \mathcal{Y}^*$ .
  4. For  $t$  from  $n$  to 2:
    - Let  $\delta_{\pi(t-1)} = \delta_{\pi(t)} - \frac{\alpha_{\pi(t)}}{\beta^t} - \sum_{\tau=1}^{t-2} \mu_{\pi(\tau)}$ .
    - Let  $\mathcal{Y}_{\pi(t-1)}$  be such that  $s(\mathcal{Y}_{\pi(t-1)}) = \delta_{\pi(t-1)}$ .
  5. Output  $\omega = (\pi, (\mathcal{Y}_{\pi(1)}, \dots, \mathcal{Y}_{\pi(n)}))$ .
- 

**PROOF.** The fact that the algorithm runs in polynomial time is immediate, since:

- additions, comparisons, and finding minimum weight matchings in a graph [Edm65] can all be done in (randomized) polynomial time; and
- the Efficient Output Divisibility Condition implies that computing a distribution  $\mathcal{Y}_{\pi(t-1)}$  such that  $s(\mathcal{Y}_{\pi(t-1)}) = \delta_{\pi(t-1)}$  is efficient.

Recall that  $(\vec{x}, \vec{\lambda})$  supports a collaborative equilibrium if and only if there exists a permutation  $\pi$  such that

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n (n - t)\mu_{\pi(t)} \leq s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0).$$

Note that our algorithm constructs a complete bipartite graph  $G = (L \cup R, E)$  where the weight on every edge is  $w(i, t) = \frac{\alpha_i}{\beta^t} + (n - t)\mu_i$ . A matching  $M$  on this graph induces a permutation  $\pi$  where, for every  $t \in R$ , we have  $\pi(t) = i$  such that  $(i, t) \in M$ . The weight of such a matching is

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n (n - t)\mu_{\pi(t)}.$$

Thus,  $(\vec{x}, \vec{\lambda})$  supports a collaborative equilibrium if and only if the maximum-weight matching in  $G$  has weight less than or equal to  $w^* \stackrel{\text{def}}{=} s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0)$ . Note that when the weight of the maximum-matching is greater than  $w^*$ , our algorithm outputs NONE, indicating that an equilibrium is not supported by the inputs.

Finally, when the weight of the maximum matching is less than or equal to  $w^*$ , the algorithm outputs a pair  $(\pi, \vec{\mathcal{Y}})$  which (by construction) satisfies  $s(\mathcal{Y}_{\pi(t-1)}) - s(\mathcal{Y}_{\pi(t)}) = \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{\tau=1}^{t-2} \mu_{\pi(\tau)}$  for all  $t \in [n]$ , so the sufficient conditions for  $(\pi, \vec{\mathcal{Y}})$  to be a collaborative equilibrium are satisfied.  $\square$

## 2.7 General NP-completeness

One may wonder if we can get an efficient mechanism for learning vectors which are not  $n$ -dimensional. We show that this is unlikely, since finding a collaborative equilibrium is NP-complete even under a weak generalization of  $n$ -dimensional learning vectors.

**DEFINITION 8.** *We say that a learning vector  $\lambda \in \Lambda$  is  $n^2$ -dimensional if there exists a non-negative matrix  $(\mu_{i,j})_{(i,j) \in [n] \times [n]}$*

such that  $\lambda_{\pi, \pi(t)} = \sum_{\tau=1}^{t-1} \mu_{\pi(t), \pi(\tau)}$ . We denote by  $\Lambda_2 \subset \Lambda$  the set of all  $n^2$ -dimensional learning vectors.

When  $\lambda$  is an  $n^2$ -dimensional learning vector, the amount that player  $\pi(t)$  learns from  $\pi(\tau)$ 's output is bounded above by  $\mu_{\pi(t), \pi(\tau)}$ . Thus, the total amount that player  $\pi(t)$  learns from all prior outputs is  $\sum_{\tau=1}^{t-1} \mu_{\pi(t), \pi(\tau)}$ . The corresponding necessary condition for a collaborative equilibrium to be supported by some  $(\vec{x}, \vec{\lambda})$  is that there is a permutation  $\pi$  such that

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n \sum_{s>t} \mu_{\pi(s), \pi(t)} \leq s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0).$$

We show that even checking whether this condition holds is NP-complete.

**THEOREM 3.** *Given model parameters CK, input datasets  $(x_1, \dots, x_n) \in X$ , and a  $n^2$ -dimensional learning bound vector  $(\mu_{i,j})_{(i,j) \in [n] \times [n]}$ , it is NP-complete to decide whether there exists  $\pi$  such that*

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n \sum_{s>t} \mu_{\pi(s), \pi(t)} \leq s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0).$$

**PROOF.** It is clear that the problem is in NP, since given a permutation  $\pi$ , the left-hand side can be efficiently computed and compared to the right-hand side of the inequality.

To show that the problem is NP-hard, we reduce it to the *minimum weighted feedback arc set problem*. The unweighted version of this problem was shown to be NP-complete by Karp [Kar72], and the weighted version is also NP-complete [ENSS95].

#### Minimum-Weight-Feedback-Arc-Set

INPUTS: A graph  $G = (V, E)$  and a weight function  $w : E \rightarrow \mathbb{R}_{\geq 0}$ , a threshold  $\gamma \in \mathbb{R}_{\geq 0}$ .

OUTPUT: Whether or not there exists a set  $S \subset E$  of edges which intersects every cycle of  $G$  and has weight less than  $\gamma$ .

All that we need to show is that, given a graph  $G$ , a set  $S$  of edges is a feedback arc set if and only if there exists a permutation  $\pi$  of the vertices of  $V$  such that  $S = \{(\pi(t), \pi(s)) \in E : s < t\}$ .

To see this, note that if  $\pi$  is a permutation and  $S = \{(\pi(t), \pi(s)) \in E : s < t\}$  then the set  $S$  intersects every cycle of  $G$ . This is because, if

$$C = \{(\pi(i_1), \pi(i_2)), (\pi(i_2), \pi(i_3)), \dots, (\pi(i_k), \pi(i_1))\}$$

is a cycle in  $G$ , then there must exist  $s, t$  such that  $s < t$  and  $(\pi(t), \pi(s)) \in C$ , so  $S$  intersects  $C$ . Thus,  $S$  is a feedback arc set.

Conversely, if  $S$  is a feedback arc set, then  $G' = (V, E - S)$  is a directed acyclic graph, and we can induce an ordering  $\pi$  on  $V$  following topological sort. Any edge  $(\pi(t), \pi(s)) \in E - S$  must satisfy  $t < s$ . Thus, any edge  $(\pi(t), \pi(s))$  where  $s < t$  must be in  $S$ . Thus, given  $\pi$  from the topological sort, we must have  $S \supset \{(\pi(t), \pi(s)) \in E : s < t\}$ . Since weights are non-negative, the minimal feedback arc set  $S^*$  will correspond to a permutation  $\pi^*$  such that  $S^* = \{(\pi^*(t), \pi^*(s)) \in E : s < t\}$ .

We show how to reduce Minimum-Weight-Feedback-Arc-Set to our problem. Given  $G = (V, E)$ ,  $w : E \rightarrow \mathbb{R}_{\geq 0}$  and  $t \in$

$\mathbb{R}_{\geq 0}$ , let  $s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0) = \gamma$  and let  $\mu_{i,j} = w(i, j)$  if  $(i, j) \in E$  and  $\mu_{i,j} = 0$  otherwise. Suppose there exists a permutation  $\pi$  such that

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n \sum_{s>t} \mu_{\pi(s), \pi(t)} \leq s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0).$$

Then, since the  $\alpha_i$  and  $\beta$  are positive,

$$\sum_{t=1}^n \sum_{s>t} \mu_{\pi(s), \pi(t)} \leq s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0).$$

Plugging in our choices of  $\mu_{i,j}$  and  $s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0)$ , this becomes

$$\sum_{(\pi(s), \pi(t)) \in E : s>t} w(\pi(s), \pi(t)) \leq \gamma.$$

Since the set  $S = \{(\pi(s), \pi(t)) \in E : s > t\}$  is a feedback arc set, we have that there exists a feedback arc set with weight less than  $\gamma$ .

Conversely, assume no such permutation  $\pi$  exists. That is,

$$\sum_{(\pi(s), \pi(t)) : s>t} \mu_{\pi(s), \pi(t)} > s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0)$$

for all permutations  $\pi$ . Note that whether  $s$  comes before  $t$  or vice-versa does not matter, since this inequality holds for all permutations. Thus, we can also write

$$\sum_{(\pi(s), \pi(t)) : s<t} \mu_{\pi(s), \pi(t)} > s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0)$$

for all permutations  $\pi$ . From the argument above, the minimum weight feedback arc set  $S^*$  induces a permutation  $\pi^*$  such that  $S^* = \{(\pi^*(t), \pi^*(s)) \in E : s < t\}$ . The weight of  $S^*$  is

$$\begin{aligned} \sum_{(\pi^*(s), \pi^*(t)) \in E : s<t} w(\pi^*(s), \pi^*(t)) &= \\ \sum_{(\pi^*(s), \pi^*(t)) : s<t} \mu_{\pi^*(s), \pi^*(t)} &> s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0) = \gamma. \end{aligned}$$

Thus, there does not exist a feedback arc set with weight less than or equal to  $\gamma$ .

We conclude that if we can efficiently check whether

$$\sum_{t=1}^n \frac{\alpha_{\pi(t)}}{\beta^t} + \sum_{t=1}^n \sum_{s>t} \mu_{\pi(s), \pi(t)} \leq s(\{\hat{y}|x_1, \dots, x_n\}) - s(\mathcal{Y}_0),$$

then we can efficiently check whether there exists a feedback arc set  $S$  with weight less than  $\gamma$ . Thus, the feedback arc set problem reduces to ours, and our problem is NP-complete.  $\square$

We have shown that in our model of scientific collaboration, it can indeed be very beneficial to *all parties involved* to collaborate under certain ordering functions, and such beneficial collaboration outcomes can be efficiently computed under certain realistic conditions (but probably not in the general case).

### 3. ORDERED MPC

We introduce formal definitions of ordered MPC and associated notions of fairness and ordered output delivery, and give protocols that realize these notions. Our definitions

build upon the standard security notion<sup>16</sup> for traditional MPC, which is described formally in Appendix ??.

### Notation.

For a finite set  $A$ , we will write  $a \leftarrow A$  to denote that  $a$  is drawn uniformly at random from  $A$ . For  $n \in \mathbb{N}$ ,  $[n]$  denotes the set  $\{1, 2, \dots, n\}$ . The operation  $\oplus$  stands for exclusive-or. The relation  $\approx$  denotes computational indistinguishability.  $\text{negl}(n)$  denotes a negligible function in  $n$ , and  $\text{poly}(n)$  denotes a polynomial in  $n$ .  $\circ$  denotes function composition, and for a function  $f$ , we write  $f^t$  to denote  $\underbrace{f \circ f \circ \dots \circ f}_t$ .

Throughout this work, we consider computationally bounded (rushing) adversaries in a synchronous complete network, and we assume the players are honest-but-curious, since any protocol secure in the presence of honest-but-curious players can be transformed into a protocol secure against malicious players [GMW87].

## 3.1 Definitions

Let  $f$  be an arbitrary  $n$ -ary function and  $p$  be an  $n$ -ary function that outputs permutation  $[n] \rightarrow [n]$ . An ordered MPC protocol is executed by  $n$  parties, where each party  $i \in [n]$  has a private input  $x_i \in \{0, 1\}^*$ , who wish to securely compute  $f(x_1, \dots, x_n) = (y_1, \dots, y_n) \in (\{0, 1\}^*)^n$  where  $y_i$  is the output of party  $i$ . Moreover, the parties are to receive their outputs in a particular *ordering* dictated by  $p(x_1, \dots, x_n) = \pi \in ([n] \rightarrow [n])$ . That is, for all  $i < j$ , party  $\pi(i)$  must receive his output *before* party  $\pi(j)$  receives her output. Note that the output ordering  $\pi$  is *data-dependent*, as  $p$  is a function of the parties' inputs.

Following [GMW87], the security of ordered MPC with respect to a functionality  $f$  and permutation function  $p$  is defined by comparing the execution of a protocol to an ideal process  $\mathcal{F}_{\text{Ordered-MPC}}$  where the outputs and ordering are computed by a trusted party who sees all the inputs. An ordered MPC protocol  $F$  is considered to be secure if for any real-world adversary  $\mathcal{A}$  attacking the real protocol  $F$ , there exists an ideal adversary  $\mathcal{S}$  in the ideal process whose outputs (views) are indistinguishable from those of  $\mathcal{A}$ . Note that this implies that no player learns more information about the other players' inputs than can be learned from his own input and output, *and his own position in the output delivery order*. The latter condition is important because the output ordering depends on parties' private inputs, and thus we require that the protocol reveals as little information as possible about the ordering.

### Many rather than one view.

In the ordered MPC setting, the ideal adversary  $\mathcal{S}$  and the real-world adversary  $\mathcal{A}$  each output a view after each output phase. This is in contrast to standard MPC, where the adversaries simply output one view at the end of the protocol execution.

IDEAL FUNCTIONALITY  $\mathcal{F}_{\text{Ordered-MPC}}$

In the ideal model, a trusted third party  $T$  is given the inputs, computes the functions  $f, p$  on the inputs, and outputs to each player  $i$  his output  $y_i$  in the order prescribed by the ordering function. In addition, we model an ideal process adversary  $\mathcal{S}$  who attacks the protocol by corrupting players in the ideal setting.

*Public parameters.*  $\kappa \in \mathbb{N}$ , the security parameter;  $n \in \mathbb{N}$ , the number of parties;  $f : (\{0, 1\}^*)^n \rightarrow (\{0, 1\}^*)^n$ , the function to be computed; and  $p : (\{0, 1\}^*)^n \rightarrow ([n] \rightarrow [n])$ , the ordering function.

*Private parameters.* Each player  $i \in [n]$  holds a private input  $x_i \in \{0, 1\}^*$ .

1. INPUT. Each player  $i$  sends his input  $x_i$  to  $T$ .
2. COMPUTATION.  $T$  computes  $(y_1, \dots, y_n) = f(x_1, \dots, x_n)$  and  $\pi = p(x_1, \dots, x_n)$ .
3. OUTPUT. The output proceeds in  $n$  sequential output rounds. At the start of the  $j^{\text{th}}$  round,  $T$  sends the output value  $\text{out}_{i,j}$  to each party  $i$ , where  $\text{out}_{j,j} = y_{\pi(j)}$  and  $\text{out}_{i,j} = \perp$  for all  $i \neq j$ . When party  $\pi(j)$  receives his output, he responds to  $T$  with the message *ack*. (The players who receive  $\perp$  are not expected to respond.) Upon receipt of the *ack*,  $T$  proceeds to the  $(j+1)^{\text{th}}$  round – or, if  $j = n$ , then the protocol terminates.
4. OUTPUT OF VIEWS. At each output round, after receiving his message from  $T$ , each party produces an output, as follows. Each uncorrupted party  $i$  outputs  $y_i$  if he has already received his output, or  $\perp$  if he has not. Each corrupted party outputs  $\perp$ . Additionally, the adversary  $\mathcal{S}$  outputs an arbitrary function of the information that he has learned during the execution of the ideal protocol.

Let the output of party  $i$  in the  $j^{\text{th}}$  round be denoted by  $\mathcal{V}_{i,j}$ , and let the view outputted by  $\mathcal{S}$  in the  $j^{\text{th}}$  round be denoted by  $\mathcal{V}_{\mathcal{S},j}$ . Let  $\mathcal{V}_{\text{Ordered-MPC}}^{\text{ideal}}$  denote the collection of all views for all output rounds:

$$(\mathcal{V}_{\mathcal{S},1}, \mathcal{V}_{1,1}, \dots, \mathcal{V}_{n,1}), \dots, (\mathcal{V}_{\mathcal{S},n}, \mathcal{V}_{1,n}, \dots, \mathcal{V}_{n,n}).$$

(If the protocol is terminated early, then views for rounds which have not yet been started are taken to be  $\perp$ .)

**DEFINITION 9 (SECURITY).** A multi-party protocol  $F$  is said to securely realize  $\mathcal{F}_{\text{Ordered-MPC}}$ , if the following conditions hold.

1. The protocol description specifies  $n$  check-points  $C_1, \dots, C_n$  corresponding to events during the execution of the protocol.
2. Take any PPT adversary  $\mathcal{A}$  who corrupts a subset of players  $S \subset [n]$ , and let  $V_{\mathcal{A},j}$  be the result of an arbitrary function  $A$  applies to his view after each check-point  $C_j$ . Let

$$V_{\mathcal{A}}^{\text{real}} = ((V_{\mathcal{A},1}, V_{1,1}, \dots, V_{n,1}), \dots, (V_{\mathcal{A},n}, V_{1,n}, \dots, V_{n,n}))$$

be the tuple consisting of the adversary  $\mathcal{A}$ 's outputted views along with the outputs of the real-world parties as specified in the ideal functionality description. Then there is a PPT ideal adversary  $\mathcal{S}$  which, attacking  $\mathcal{F}_{\text{Ordered-MPC}}$  by corrupting the same subset  $S$  of players, can output

<sup>16</sup>Note that throughout this work, we use “stand-alone” security notions rather than “universally composable” ones.

views  $\mathcal{V}_{S,j}$  such that for each  $j \in [n]$ , it holds that  $\mathcal{V}_{S,j} \stackrel{c}{\approx} V_{A,j}$ .

In the context of ordered MPC, the standard guaranteed output delivery notion is insufficient. Instead, we define *ordered output delivery*, which requires in addition that all parties receive their outputs in the order prescribed by  $p$ .

**DEFINITION 10 (ORDERED OUTPUT DELIVERY).** *An ordered MPC protocol satisfies ordered output delivery if for any inputs  $x_1, \dots, x_n$ , functionality  $f$ , and ordering function  $p$ , it holds that all parties receive their outputs before protocol termination, and moreover, if  $\pi(i) < \pi(j)$ , then party  $i$  receives his output before party  $j$  receives hers, where  $\pi = p(x_1, \dots, x_n)$ .*

We also define a natural relaxation of the fairness requirement for ordered MPC, called *prefix-fairness*. Although it is known that fairness is impossible for general functionalities in the presence of a dishonest majority, we show in the next subsection that prefix-fairness can be achieved even when a majority of parties are corrupt. We emphasize that this notion relaxes *only* the fairness requirement: that is, prefix-fair protocols satisfy full privacy (and correctness) guarantees.

**DEFINITION 11 (PREFIX-FAIRNESS).** *An ordered MPC protocol is prefix-fair if for any inputs  $x_1, \dots, x_n$ , it holds that the set of parties who have received their outputs at the time of protocol termination (or abortion) is a prefix of  $(\pi(1), \dots, \pi(n))$ , where  $\pi = p(x_1, \dots, x_n)$  is the permutation induced by the inputs.*

Prefix-fairness can be useful, for example, in settings where it is more important for one party to receive the output than the other; or where there is some prior knowledge about the trustworthiness of each party (so that more trustworthy parties may receive their outputs first).

### 3.2 Construction

Ordered MPC is achievable by using standard protocols for general MPC, as described in Protocol 1 below. The protocol has  $n$  sequential output phases, so that the  $n$  outputs can be issued in order. A subtle point is that because the ordering is a function of the input data, knowledge of the ordering may reveal information about the input data. Thus, we have to “mask” the output values such that each party only learns the minimal possible amount of information about the ordering: namely, his own position in the ordering.

#### PROTOCOL 1: ORDERED MPC

*Public parameters.*  $\kappa \in \mathbb{N}$ , the security parameter;  $n \in \mathbb{N}$ , the number of parties;  $k \in \mathbb{N}$ , an upper bound on the number of corrupt parties;  $f : (\{0, 1\}^*)^n \rightarrow (\{0, 1\}^*)^n$ , the function to be computed; and  $p : (\{0, 1\}^*)^n \rightarrow [n]$ , the ordering function.

1. **Computing shares of  $(\pi, \mathbf{y})$ :** Using any general secure MPC protocol (such as [GMW87]) on inputs  $x_1, \dots, x_n$ , jointly compute a  $k$ -out-of- $n$  secret-sharing<sup>17</sup> of  $(\pi, \mathbf{y})$  where  $\mathbf{y} = (y_1, \dots, y_n) =$

$f(x_1, \dots, x_n)$  and  $\pi = p(x_1, \dots, x_n)$  is a permutation of  $[n]$ . At the end of this step, each player possesses a share of the outputs  $\mathbf{y} = (y_1, \dots, y_n)$  and of the permutation  $\pi$ .

2. **Outputting  $y_1, \dots, y_n$  in  $n$  phases:** In the  $i^{\text{th}}$  output phase, player  $\pi^{-1}(i)$  will learn his output. In phase  $i$  the parties run a new instance of a general secure MPC protocol such that:

- Player  $j$ 's inputs to the protocol are: the shares of  $\mathbf{y}$  and  $\pi$  that he got in step 1, and a random string  $r_{i,j}$ .
- The functionality computed is:

```
for  $j$  from 1 to  $n$ : if  $\pi(j) = i$ 
then  $z_{i,j} := y_j \oplus r_{i,j}$  else  $z_{i,j} = \perp \oplus r_{i,j}$ .
output  $z_i = (z_{i,1}, \dots, z_{i,n})$ .
```

where  $\perp$  is a special string that lies outside the output domain.

- To recover his output, each player  $j$  computes  $y'_{i,j} = z_{i,j} \oplus r_{i,j}$  for all  $i$ . By construction, there is exactly one  $i \in [n]$  for which  $y'_{i,j} \neq \perp$ , and that is equal to the output value  $y_j$  for player  $j$ .

*Check-points.* There are  $n$  check-points. For  $i \in [n]$ , the check-point  $C_i$  is at the end of the  $i^{\text{th}}$  output phase, when  $z_i$  is learned by all players.

*In case of abort.* When running the protocol for the honest majority setting, the honest players continue until the end of the protocol regardless of other players' behavior. When running the protocol for dishonest majority, if any party aborts in an output phase<sup>18</sup>, then the honest players do not continue to the next phase.

In proving the security of Protocol 1, we refer to the security of modular composition of general protocols shown by [Can00], Theorem 5.

**THEOREM 4.** *Protocol 1 securely realizes  $\mathcal{F}_{\text{Ordered-MPC}}$ .*

**PROOF.** Let  $\rho_0$  denote the general MPC protocol execution in step 1, and let  $\rho_i$  be the general MPC protocol execution in phase  $i$  of step 2, for  $i \in [n]$ . For  $j \in [n]$ , let the protocol  $\pi_j$  be the concatenation of the protocols  $\rho_0, \dots, \rho_j$ . To prove security at each check-point, it is sufficient to prove that  $\pi_j$  satisfies security for all  $j \in [n]$ : in other words, that the view outputted by any adversary in the real protocol execution at check-point  $j$  can be simulated in the ideal execution. Finally, for all  $j \in [n]$ , the security of  $\pi_j$  follows directly from the security of modular composition of general protocols ([Can00], Theorem 5).  $\square$

<sup>17</sup>The standard definition of a secret-sharing scheme can be found in Appendix ??.

<sup>18</sup>Each output phase consists of an execution of the underlying general MPC protocol. If a party aborts at any time during (and before the end of) the execution of the underlying general MPC protocol, this fact will be detected by all honest parties by the end of the phase.

**THEOREM 5.** *In the case of honest majority, Protocol 1 achieves fairness. In the dishonest majority setting, prefix-fairness is achieved.*

**PROOF.** Fairness holds in the honest majority case, since the honest players complete all output phases, and the shares that the honest players hold are sufficient to reconstruct each output  $y_i$  (recall that the secret-sharing threshold  $k$  is  $\lceil n/2 \rceil$  in the honest majority case). In the dishonest majority setting, prefix-fairness holds since for all  $i \in [n]$ , all  $n$  shares are required in order to reconstruct the output  $y_{\pi(i)}$  in output phase  $i$ , and

- if the corrupt parties do not abort during the  $i^{\text{th}}$  output phase, then by the security of Protocol 1, the output  $y_{\pi(i)}$  associated with the  $i^{\text{th}}$  output phase is delivered correctly to party  $i$ ;
- if the corrupt parties abort during the  $i^{\text{th}}$  output phase, then no outputs  $y_{\pi(j)}$  for  $j > i$  will be learned by any player, since the honest parties will not execute subsequent output phases.

□

## 4. TIMED-DELAY MPC

In this section, we implement *time delays* between different players receiving their outputs. The model is exactly as before, with  $n$  players wishing to compute a function  $f(x_1, \dots, x_n)$  in an ordering prescribed by  $p(x_1, \dots, x_n)$  – except that now, there is an additional requirement of a delay after each player receives his output and before the next player receives her output. To realize the timed-delay MPC functionality, we make use of time-lock and time-line puzzles, which are introduced in Sections 4.3.1 and 4.4.

### 4.1 Ideal functionality with time delays

We measure time delay in units of computation, rather than seconds of a clock: that is, rather than making any assumption about global clocks (or synchrony of local clocks)<sup>19</sup>, we measure time by the *evaluations of a particular function* (on random inputs), which we call the *clock function*.

#### IDEAL FUNCTIONALITY $\mathcal{F}_{\text{Timed-Delay-MPC}}$

In the ideal model, a trusted third party  $T$  is given the inputs, computes the functions  $f, p$  on the inputs, and outputs to each player  $i$  his output  $y_i$  in the order prescribed by the

<sup>19</sup>A particular issue that arises when considering a clock-based definition is that it is not clear that we can reasonably assume or prove that clocks are in synchrony between the real and ideal world – but this seems necessary in order to prove security by simulation in the ideal functionality.

We remark that if one is happy to assume the existence of a global clock (or synchrony of local clocks), then there are other ways to implement timed-delay MPC which sidestep many of the issues inherent in the arguably more realistic model where clocks may not be perfectly synchronized between different (adversarial) parties. One example is the “Bitcoin model” where the assumption is that the Bitcoin block-chain can serve as a global clock: in this model, existing protocols such as [BK14] implement some time-delays in MPC, and it seems likely that such protocols can be adapted to achieve our notion of timed-delay MPC.

ordering function. Moreover,  $T$  imposes delays between the issuance of one party’s output and the next. In addition, we model an ideal process adversary  $\mathcal{S}$  who attacks the protocol by corrupting players in the ideal setting.

**Public parameters.**  $\kappa \in \mathbb{N}$ , the security parameter;  $n \in \mathbb{N}$ , the number of parties;  $f : (\{0, 1\}^*)^n \rightarrow (\{0, 1\}^*)^n$ , the function to be computed;  $p : (\{0, 1\}^*)^n \rightarrow [n] \rightarrow [n]$ , the ordering function; and  $G = G(\kappa) \in \mathbb{N}$ , the number of time-steps between the issuance of one party’s output and the next.

**Private parameters.** Each player  $i \in [n]$  holds a private input  $x_i \in \{0, 1\}^*$ .

1. **INPUT.** Each player  $i$  sends his input  $x_i$  to  $T$ . If, instead of sending his input, any player sends the message quit, then the computation is aborted.
2. **COMPUTATION.**  $T$  computes  $(y_1, \dots, y_n) = f(x_1, \dots, x_n)$  and  $\pi = p(x_1, \dots, x_n)$ .
3. **OUTPUT.** The output proceeds in  $n$  sequential output phases. At each phase  $j$ ,  $T$  waits for  $G$  time-steps, then sends the  $j^{\text{th}}$  output,  $y_{\pi(j)}$ , to party  $\pi(j)$ .
4. **OUTPUT OF VIEWS.** At the end of each output phase, each party produces an output as follows. Each uncorrupted party  $i$  outputs  $y_i$  as his view if he has already received his output, or  $\perp$  if he has not. Each corrupted party outputs  $\perp$ . Additionally, the adversary  $\mathcal{S}$  outputs an arbitrary function of the information that he has learned during the execution of the ideal protocol, after each check-point.

Let the output of party  $i$  in the  $j^{\text{th}}$  round be denoted by  $\mathcal{V}_{i,j}$ , and let the view outputted by  $\mathcal{S}$  in the  $j^{\text{th}}$  round be denoted by  $\mathcal{V}_{\mathcal{S},j}$ . Let  $\mathcal{V}_{\text{Timed-Delay-MPC}}^{\text{ideal}}$  denote the collection of all views for all output phases:

$$\mathcal{V}_{\text{Timed-Delay-MPC}}^{\text{ideal}} = ((\mathcal{V}_{\mathcal{S},1}, \mathcal{V}_{1,1}, \dots, \mathcal{V}_{n,1}), \dots, (\mathcal{V}_{\mathcal{S},n}, \mathcal{V}_{1,n}, \dots, \mathcal{V}_{n,n})).$$

For an algorithm  $\mathcal{A}$ , let the run-time<sup>20</sup> of  $\mathcal{A}$  on input  $\text{inp}$  be denoted by  $\text{time}_{\mathcal{A}}(\text{inp})$ . If  $\mathcal{A}$  is probabilistic, the run-time will be a distribution over the random coins of  $\mathcal{A}$ . Note that the exact run-time of an algorithm will depend on the underlying computational model in which the algorithm is run. In this work, all algorithms are assumed to be running in the same underlying computational model, and our definitions and results hold regardless of the specific computational model employed.

**DEFINITION 12 (SECURITY).** *A multi-party protocol  $F$  (with parameters  $\kappa, n, f, p, G$ ) is said to securely realize  $\mathcal{F}_{\text{Timed-Delay-MPC}}$ , if the following conditions hold.*

1. *The protocol description specifies  $n$  check-points  $C_1, \dots, C_n$  corresponding to events during the execution of the protocol.*
2. *There exists a “clock function”  $g$  such that between any two consecutive checkpoints  $C_i, C_{i+1}$  during an execution of  $F$ , any one of the parties (in the real world) must be able to locally run  $\Omega(G)$  sequential evaluations*

<sup>19</sup>The use of checkpoints is introduced to capture the views of players and the adversary at intermediate points in protocol execution.

<sup>20</sup>Run-time is, naturally, measured in “CPU time” (i.e. the number of instructions executed in the underlying computational model) as opposed to real-world “clock time”.

of  $g$  on random inputs.  $g$  may also be a protocol (involving  $n' \leq n$  parties) rather than a function, in which case we instead require that any subset consisting of  $n'$  parties must be able to run  $\Omega(G)$  sequential executions of  $g$  (on random inputs) over the communication network being used for the main multi-party protocol  $F$ . Then, we say that  $F$  is “clocked by  $g$ ”.

3. Take any PPT adversary  $\mathcal{A}$  attacking the protocol  $F$  by corrupting a subset of players  $S \subset [n]$ , which outputs an arbitrary function  $V_{\mathcal{A},j}$  of the information that it has learned in the protocol execution after each check-point  $C_j$ . Let

$$V_{\mathcal{A}}^{\text{real}} = ((V_{\mathcal{A},1}, V_{1,1}, \dots, V_{n,1}), \dots, (V_{\mathcal{A},n}, V_{1,n}, \dots, V_{n,n}))$$

be the tuple consisting of the adversary  $\mathcal{A}$ 's outputted views along with the views of the real-world parties as specified in the ideal functionality description. Then there is a PPT ideal adversary  $\mathcal{S}$  which, attacking  $\mathcal{F}_{\text{Timed-Delay-MPC}}$  by corrupting the same subset  $S$  of players, can output views  $\mathcal{V}_{\mathcal{S},1}, \dots, \mathcal{V}_{\mathcal{S},n}$  (at check-points  $C_1, \dots, C_n$  respectively) such that for each  $j \in [n]$ , it holds that

$$\begin{aligned} & |\Pr[D(\mathcal{V}_{\mathcal{S},j}, \mathcal{V}_{1,j}, \dots, \mathcal{V}_{n,j}) = 1] \\ & - \Pr[D(V_{\mathcal{A},j}, V_{1,j}, \dots, V_{n,j}) = 1]| \\ & \leq \text{negl}(\kappa), \end{aligned}$$

for any distinguisher  $D$  such that

$$\Pr_{\vec{v} \leftarrow \mathcal{V}}[\text{time}_D(\vec{v}) \leq j \cdot \text{time}_G()] = 1/\text{poly}(\kappa),$$

when  $\mathcal{V}$  is the distribution of views outputted by  $\mathcal{A}$  or  $\mathcal{S}$  (that is, for  $\mathcal{V} \in \{(\mathcal{V}_{\mathcal{S},j}, \mathcal{V}_{1,j}, \dots, \mathcal{V}_{n,j}), (V_{\mathcal{A},j}, V_{1,j}, \dots, V_{n,j})\}$ ), and  $G$  is the algorithm that computes the function  $g$  sequentially on  $G$  random inputs.

## 4.2 Realizing timed-delay MPC with dummy rounds

A simple protocol for securely realizing timed-delay MPC is to implement delays by running  $G$  “dummy rounds” of communication in between issuing outputs to different players.

### PROTOCOL 2: TIMED-DELAY MPC WITH DUMMY ROUNDS

*Public parameters.*  $\kappa \in \mathbb{N}$ , the security parameter;  $n \in \mathbb{N}$ , the number of parties;  $f : (\{0, 1\}^*)^n \rightarrow (\{0, 1\}^*)^n$ , the function to be computed;  $p : (\{0, 1\}^*)^n \rightarrow ([n] \rightarrow [n])$ , the ordering function; and  $G = \text{poly}(\kappa)$ , the number of time-steps between the issuance of one party's output and the next.

1. **Computing shares of  $(\pi, \mathbf{y})$ :** Using any general secure MPC protocol (such as [GMW87]), jointly compute an  $k$ -out-of- $n$  secret-sharing<sup>21</sup> of  $(\pi, \mathbf{y})$  where  $\mathbf{y} = (y_1, \dots, y_n) = f(x_1, \dots, x_n)$  and permutation  $\pi = p(x_1, \dots, x_n)$  on the players' inputs. At the end of this step, each player possesses a share of the outputs  $\mathbf{y} = (y_1, \dots, y_n)$  and of the permutation  $\pi$ .

2. **Outputting  $y_1, \dots, y_n$  in  $n$  phases:** The outputs will occur in  $n$  phases: in the  $i^{\text{th}}$  phase, player  $\pi^{-1}(i)$  will learn his output. In each phase, the players first run  $G$  “dummy rounds” of communication. A dummy round is a “mini-protocol” defined as follows (let this mini-protocol be denoted by  $g_{\text{dum}}$ ):
  - each player initially sends the message **challenge** to every other player;
  - each player responds to each **challenge** he receives with a message **response**.

In each phase, after the dummy rounds have been completed, the parties will run a new instance of a general secure MPC protocol. In phase  $i$ :

- Player  $j$ 's inputs to the protocol are: the shares of  $\mathbf{y}$  and  $\pi$  that he got in step 1, and a fresh random string  $r_{i,j}$ .
- The functionality computed in each phase  $i \in [n]$  is:
 

```

      for  $j$  from 1 to  $n$ : if  $\pi(j) = i$ 
      then  $z_{i,j} := y_j \oplus r_{i,j}$  else  $z_{i,j} = \perp \oplus r_{i,j}$ 
      output  $z_i = (z_{i,1}, \dots, z_{i,n})$ .
      
```

 where  $\perp$  is a special string that lies outside the output domain.
- To recover his output, each player  $j$  computes  $y'_{i,j} = z_{i,j} \oplus r_{i,j}$  for all  $i$ . By construction, there is exactly one  $i \in [n]$  for which  $y'_{i,j} \neq \perp$ , and that is equal to the output value  $y_j$  for player  $j$ .

*Check-points.* There are  $n$  check-points. For  $i \in [n]$ , the check-point  $C_i$  is at the end of the  $i^{\text{th}}$  output phase, when  $z_i$  is learned by all players.

*In case of abort.* When running the protocol for the honest majority setting, the honest players continue until the end of the protocol regardless of other players' behavior. When running the protocol for dishonest majority, if any party aborts in an output phase<sup>22</sup>, then the honest players do not continue to the next phase.

**THEOREM 6.** *In the presence of honest majority, Protocol 2 securely realizes  $\mathcal{F}_{\text{Timed-Delay-MPC}}$  clocked by  $g_{\text{dum}}$ .*

<sup>21</sup>For the honest majority setting, we set  $k = \lceil n/2 \rceil$ . For the dishonest majority setting,  $k = n$ .

<sup>22</sup>Each output phase consists of an execution of the underlying general MPC protocol preceded by  $G$  dummy rounds. If a party aborts before the completion of the  $G$  dummy rounds, this fact will be detected by all parties in the dummy round in which the abort happens, because every party is supposed to communicate with every other party in each dummy round. If a party aborts at any time during (and before the end of) the execution of the underlying general MPC protocol, this fact will be detected by all honest parties by the end of the phase.



PROOF. Let  $\mathcal{A}$  be any PPT adversary attacking Protocol 2 by corrupting a subset of players  $S \subset [n]$ , and let

$$V_{\mathcal{A}}^{\text{real}} = ((V_{\mathcal{A},1}, V_{1,1}, \dots, V_{n,1}), \dots, (V_{\mathcal{A},n}, V_{1,n}, \dots, V_{n,n}))$$

be the tuple consisting of the adversary  $\mathcal{A}$ 's outputted views along with the views of the real-world parties (as specified in the description of  $\mathcal{F}_{\text{Timed-Delay-MPC}}$ ). In order to show that condition 3 of the security definition (Definition 12) holds, we need to show that there is a PPT ideal adversary  $\mathcal{S}$  which, given access to  $\mathcal{F}_{\text{Timed-Delay-MPC}}$  and corrupting the same subset  $S$  of players, can output views  $\mathcal{V}_{S,j}$  such that  $V_{\mathcal{A}}^{\text{real}} \stackrel{c}{\approx} \mathcal{V}_{\text{Timed-Delay-MPC}}^{\text{ideal}}$ .

Recall that the adversary's view can be any function of the inputs of the corrupt parties and the messages that the corrupt parties see during the protocol execution. In particular, it is sufficient to show that there is an ideal adversary  $\mathcal{S}$  which can output views  $\mathcal{V}_{S,j}$  which are indistinguishable from the transcript of all the messages that the corrupt parties see during the real protocol execution.

Protocol 2 consists of sequential executions of the underlying general MPC protocol and the mini-protocol  $g_{\text{dum}}$ . When the mini-protocol executions are removed from Protocol 2, the resulting protocol is identical to Protocol 1. Hence, by Theorem 4, there is an ideal adversary  $\mathcal{S}'$  which, given access to  $\mathcal{F}_{\text{Timed-Delay-MPC}}$  and corrupting the same subset  $S$  of players, can output views  $\mathcal{V}_{S',j}$  such that which are indistinguishable from the transcript of all the messages that the corrupt parties see *during the  $n+1$  executions of the underlying general MPC protocol* within Protocol 2. The only other messages that are sent in Protocol 2 are the “dummy” messages **challenge** and **response**, which are fixed messages that do not depend on the players' inputs. In fact, the transcript of an execution of  $g_{\text{dum}}$  is a deterministic sequence of **challenge** and **response**. It follows that there exists an ideal adversary  $\mathcal{S}$  which, by calling  $\mathcal{S}'$  and adding the deterministic transcript corresponding to each execution of  $g_{\text{dum}}$ , can output views  $\mathcal{V}_{S,j}$  which are indistinguishable from the transcript of all the messages that the corrupt parties see during the real execution of Protocol 2.

Finally, it remains to show that condition 2 of the security definition (Definition 12) is satisfied. The players are literally running  $g$  over the MPC network  $G$  times in between issuing outputs, so it is clear that condition 2 holds.  $\square$

One downside of the simple solution above is that it requires all (honest) parties to be online and communicating until the last player receives his output. To address this, in Section 4.3 we propose an alternative solution based on *timed-release cryptography*, at the cost of an additional assumption that all players have comparable computing speed (within a logarithmic factor).

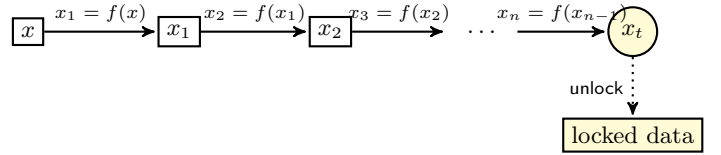
### 4.3 Realizing timed-delay MPC with time-lock puzzles

Informally, a time-lock puzzle is a primitive which allows “locking” of data, such that it will be released after a pre-specified time delay, and no earlier. Our next protocol, instead of issuing outputs to players in the clear, gives to each party his output *locked* into a time-lock puzzle; and in order to enforce the desired ordering, the delays required to unlock the puzzles are set to be an increasing sequence. We first give the definition of time-lock puzzles (in Section 4.3.1) then

describe and prove security of our time-lock-based protocol (in Section 4.3.2).

#### 4.3.1 Time-lock puzzles

The delayed release of data in MPC protocols can be closely linked to the problem of “timed-release crypto” in general, which was introduced by [May93] and constructed first by [RSW96] with their proposal of *time-lock puzzles*. We assume time-lock puzzles with a particular structure (that is present in all known implementations): namely, the passage of “time” will be measured by sequential evaluations of a function (**TimeStep**). Unlocking a  $t$ -step time-lock puzzle can be considered analogous to following a chain of  $t$  pointers, at the end of which there is a special value  $x_t$  (e.g. a decryption key) that allows retrieval of the locked data.



DEFINITION 13 (TIME-LOCK PUZZLE SCHEME). A time-lock puzzle scheme is a tuple of PPT algorithms

$$T = (\text{Lock}, \text{TimeStep}, \text{Unlock})$$

as follows:

- **Lock**( $1^\kappa, d, t$ ) takes parameters  $\kappa \in \mathbb{N}$  the security parameter,  $d \in \{0, 1\}^\ell$  the data to be locked, and  $t \in \mathbb{N}$  the number of steps needed to unlock the puzzle, and outputs a time-lock puzzle  $P = (x, t, b, a) \in \{0, 1\}^n \times \mathbb{N} \times \{0, 1\}^{n''} \times \{0, 1\}^{n'}$  where  $\ell, n, n', n'' = \text{poly}(\kappa)$ .
- **TimeStep**( $1^\kappa, x', a'$ ) takes parameters  $\kappa \in \mathbb{N}$  the security parameter, a bit-string  $x' \in \{0, 1\}^n$ , and auxiliary information  $a'$ , and outputs a bit-string  $x'' \in \{0, 1\}^n$ .
- **Unlock**( $1^\kappa, x', b'$ ) takes parameters  $\kappa \in \mathbb{N}$  the security parameter, a bit-string  $x' \in \{0, 1\}^n$ , and auxiliary information  $b' \in \{0, 1\}^{n'}$ , and outputs some data  $d' \in \{0, 1\}^\ell$ .

To unclutter notation, we will sometimes omit the initial security parameter of these functions (writing e.g. simply **Lock**( $d, t$ )). We now define some auxiliary functions. For a time-lock puzzle scheme  $T = (\text{Lock}, \text{TimeStep}, \text{Unlock})$  and  $i \in \mathbb{N}$ , let  $\text{IterateTimeStep}_i^T$  denote the following function:

$$\text{IterateTimeStep}_i^T(i, x, a) = \underbrace{\text{TimeStep}(\text{TimeStep}(\dots (\text{TimeStep}(x, a), a) \dots), a)}_i.$$

Define **CompleteUnlock** $^T$  to be the following function:

$$\text{CompleteUnlock}^T((x, t, b, a)) = \text{Unlock}(\text{IterateTimeStep}_t^T(t, x, a), b),$$

that is, the function that should be used to unlock a time-lock puzzle outputted by **Lock**.

The following definitions formalize correctness and security for time-lock puzzle schemes.

DEFINITION 14 (CORRECTNESS). A time-lock puzzle scheme  $T = (\text{Lock}, \text{TimeStep}, \text{Unlock})$  is correct if the following holds (where  $\kappa$  is the security parameter):

$$\Pr_{(x,t,b,a) \leftarrow \text{Lock}(d,t)} [\text{CompleteUnlock}^T((x,t,b,a)) \neq d] \leq \text{negl}(\kappa).$$

DEFINITION 15 (SECURITY). Let  $T = (\text{Lock}, \text{TimeStep}, \text{Unlock})$  be a time-lock puzzle scheme.  $T$  is secure if it holds that: for all  $d, d' \in \{0,1\}^\ell$ ,  $t = \text{poly}(\kappa)$ , if there exists an adversary  $\mathcal{A}$  that solves the time-lock puzzle  $\text{Lock}(d,t)$ , that is,

$$\Pr_{P \leftarrow \text{Lock}(d,t)} [\mathcal{A}(P) = d] = \varepsilon \text{ for some non-negligible } \varepsilon,$$

then for each  $j \in [t]$ , there exists an adversary  $\mathcal{A}_j$  such that

$$\Pr_{P' \leftarrow \text{Lock}(d',j)} [\mathcal{A}_j(P') = d'] \geq 1 - \text{negl}(\kappa), \text{ and}$$

$$\Pr_{\substack{P \leftarrow \text{Lock}(d,t), \\ P' \leftarrow \text{Lock}(d',j)}} [\text{time}_{\mathcal{A}}(P) \geq (t/j) \cdot \text{time}_{\mathcal{A}_j}(P') \mid \mathcal{A}(P) = d] \geq 1 - \text{negl}(\kappa).$$

### 4.3.2 Protocol based on time-lock puzzles

Because of the use of time-lock puzzles by different parties in the protocol that follows, we require an additional assumption that all players have comparable computing power (within a logarithmic factor).

#### Relative-Delay Assumption..

The difference in speed of performing computations between any two parties  $i, j \in [n]$  is at most a factor of  $B = O(\log(\kappa))$ .

#### PROTOCOL 3: TIMED-DELAY MPC WITH TIME-LOCK PUZZLES

*Public parameters.*  $\kappa \in \mathbb{N}$ , the security parameter;  $n \in \mathbb{N}$ , the number of parties;  $f : (\{0,1\}^*)^n \rightarrow (\{0,1\}^*)^n$ , the function to be computed;  $p : (\{0,1\}^*)^n \rightarrow ([n] \rightarrow [n])$ , the ordering function;  $B = O(\log(\kappa))$ , the maximum factor of difference between any two parties' computing power;  $G = \text{poly}(\kappa)$ , the number of time-steps between the issuance of one party's output and the next; and  $T = \{\text{Lock}, \text{TimeStep}, \text{Unlock}\}$  a time-lock puzzle scheme.

**INPUTS.** Each party  $i$  has input  $x_i$ .

**PROTOCOL STEPS.** Let  $(y_1, \dots, y_n) = f(x_1, \dots, x_n)$  and  $\pi = p(x_1, \dots, x_n)$ . Define  $t_1 = 1$  and  $t_{i+1} = (B \cdot G + 1) \cdot t_i$  for  $i \in [n-1]$ . Compute  $(P_1, \dots, P_n)$ , where each  $P_i = (x_i, t_{\pi(i)}, a_i, b_i)$  is a time-lock puzzle computed as

$$P_i = \text{Lock}(y_i \oplus r_i, t_{\pi(i)}),$$

where each  $r_i$  is a random string provided as input randomness by party  $i$ .

**OUTPUTS.** For each  $i \in [n]$ , the puzzle  $P_i$  is outputted to party  $i$ . The players all receive their respective outputs at the same time, then recovers his output  $y_i$  by solving his time-lock puzzle, and finally “unmasking” the result by XORing with his random input  $r_i$ .

*Check-points.* There are  $n$  check-points. For  $i \in [n]$ , the check-point  $C_i$  is the event of party  $\pi(i)$  learning his eventual output  $y_{\pi(i)}$  (i.e. when he finishes solving his time-lock puzzle).

For the following theorem, we assume that each player  $i$  uses the optimal algorithm to solve his puzzle  $P_i$  that outputs the correct answer. Without this assumption, any further protocol analysis would not make sense: there can always be a “lazy” player who willfully uses a very slow algorithm to solve his puzzle, who will as a result learn his eventual output much later in the order than he could otherwise have done. The property that we aim to achieve is that every player *could* learn his output at his assigned position in the ordering  $\pi$ , with appropriate delays before and after he learns his output.

**THEOREM 7.** Suppose that the Relative-Delay Assumption holds, and each player  $i$  uses the optimal algorithm to solve his puzzle  $P_i$  that outputs (with overwhelming probability) the correct answer. Then, Protocol 3 securely realizes  $\mathcal{F}_{\text{Timed-Delay-MPC}}$  when there is an honest majority.

**PROOF.** First, we prove that condition 2 of the security definition (Definition 12) is satisfied. Let  $\mathcal{A}_i$  denote the algorithm that party  $i$  uses to solve his time-lock puzzle, and let the time at which party  $i$  learns his answer  $y_i$  be denoted by  $\tau_i = \text{time}_{\mathcal{A}_i}(P_i)$ . By the security of the time-lock puzzles, there exists an algorithm  $\mathcal{A}'_i$  that player  $i$  could use to solve the puzzle  $\text{Lock}(0^\ell, 1)$  in time  $\tau_i/t_i$ . Moreover, by the Relative-Delay Assumption, it holds that no player can solve the puzzle  $\text{Lock}(0^\ell, 1)$  more than  $B$  times faster than another player: that is,  $\max_i(\tau_i/t_i) \leq B \cdot \min_i(\tau_i/t_i)$ . It follows that even the slowest player (call him  $i^*$ ) would be able to run  $t_i/B$  executions of  $\mathcal{A}'_{i^*}$  within time  $\tau_i$ , for any  $i$ .

Without loss of generality, assume that the ordering function  $p$  is the identity function. Consider any consecutive pair of checkpoints  $C_i, C_{i+1}$ . These checkpoints occur at times  $\tau_i$  and  $\tau_{i+1}$ , by definition. We have established that in time  $\tau_i$ , player  $i^*$  can run  $t_i/B$  executions of  $\mathcal{A}'_{i^*}$ , and in time  $\tau_{i+1}$ , he can run  $t_{i+1}/B$  executions of  $\mathcal{A}'_{i^*}$ . It follows that in between the two checkpoints (i.e. in time  $\tau_{i+1} - \tau_i$ ), he can run  $(t_{i+1} - t_i)/B$  executions of  $\mathcal{A}'_{i^*}$ . Substituting in the equation  $t_{i+1} = (B \cdot G + 1) \cdot t_i$  from the protocol definition, we get that player  $i^*$  can run  $G \cdot t_i$  executions of  $\mathcal{A}'_{i^*}$  between checkpoints  $C_i$  and  $C_{i+1}$ . Since  $t_i \geq 1$  for all  $i$ , this means that  $i^*$  can run at least  $G$  executions of  $\mathcal{A}'_{i^*}$  between *any* consecutive pair of checkpoints. Hence, condition 2 holds.

We now prove condition 3. Let  $\mathcal{G}$  be the algorithm that evaluates  $\mathcal{A}'_{i^*}$  sequentially  $G$  times on random inputs. It is sufficient to show that for any adversary  $\mathcal{A}$  attacking the protocol by corrupting a subset  $S \subset [n]$  of players, which outputs a view  $V_{\mathcal{A},j}$  at each checkpoint  $j$  which is the *transcript of all messages that it has seen so far*, there is an ideal adversary  $\mathcal{S}$  which outputs views  $\mathcal{V}_{\mathcal{S},1}, \dots, \mathcal{V}_{\mathcal{S},n}$  such that for any  $j \in [n]$ , for any distinguisher  $D$  whose run-time satisfies the conditions in Definition 12, item 3,

$$|\Pr[D(\mathcal{V}_{\mathcal{S},j}, \mathcal{V}_{1,j}, \dots, \mathcal{V}_{n,j}) = 1] - \Pr[D(V_{\mathcal{A},j}, V_{1,j}, \dots, V_{n,j}) = 1]| \leq \text{negl}(\kappa).$$

Recall that there are  $n$  sequential output stages in the ideal functionality  $\mathcal{F}_{\text{Timed-Delay-MPC}}$ . Consider an ideal adversary  $\mathcal{S}$  attacking  $\mathcal{F}_{\text{Timed-Delay-MPC}}$  by corrupt a set of parties  $S \subset [n]$ . Let  $\vec{inp}$  denote the vector of inputs and input randomness of the corrupt parties (note that these are known to  $\mathcal{S}$ ). Take any  $i \in [n]$ . In the ideal protocol execution,  $\mathcal{S}$  learns the following in the  $\pi(i)^{\text{th}}$  output stage:

- nothing, if  $i \notin S$ ; or

- the input value  $x_i$ , the input randomness  $r_i$ , and the eventual output  $y_i$  if  $i \in S$ .

Note that as a result,  $\mathcal{S}$  learns  $\pi(i)$  at output stage  $\pi(i)$ , for each  $i \in S$ . The delay values  $t_1, \dots, t_n$  are a fixed sequence of values independent of the parties' inputs, so they are known to  $\mathcal{S}$ . Thus, at each check-point  $j \in [n]$ , the ideal adversary  $\mathcal{S}$  can compute  $n$  time-lock puzzles

$$\hat{P}_{j,i} = \begin{cases} \text{Lock}(y_i \oplus r_i, t_{\pi(i)}) & \text{if } 1 \leq i \leq j \\ \text{Lock}(r_i, t_{\pi(i)}) & \text{if } j < i \leq n \end{cases}.$$

Let the ideal adversary  $\mathcal{S}$  output the following view at each check-point  $j$ :

$$\mathcal{V}_{S,j} = \mathcal{S}_j(\vec{inp}, (\hat{P}_{j,1}, \dots, \hat{P}_{j,n})),$$

where  $\mathcal{S}_j$  is the ideal adversary (for the underlying general MPC protocol) that simulates the adversary's  $j^{\text{th}}$  view  $V_{A,j}$ .

We now analyze the distribution of the puzzles  $\hat{P}_{j,i}$ . For the range  $1 \leq i \leq j$ , the puzzle  $\hat{P}_{j,i}$  is by definition identically distributed to the puzzle that is outputted to player  $i$  in the real execution of Protocol 3. Now take any  $j \in [n]$ , and let  $D$  be any distinguisher whose run-time satisfies the conditions in Definition 12, item 3. Recall that the players are assumed to solve the time-lock puzzles using the optimal algorithm. Hence, it follows from the security of the underlying time-lock puzzle scheme that for any  $i$  in the range  $j < i \leq n$ ,

$$\left| \Pr[D(P_{j,i}) = 1] - \Pr[D(\hat{P}_{j,i}) = 1] \right| \leq \text{negl}(\kappa).$$

Since we defined the outputs of  $\mathcal{S}$  to be

$$\mathcal{V}_{S,j} = \mathcal{S}_j(\vec{inp}, (\hat{P}_{j,1}, \dots, \hat{P}_{j,n}))$$

for  $j \in [n]$ , it follows that

$$\begin{aligned} |\Pr[D(\mathcal{V}_{S,j}, \mathcal{V}_{1,j}, \dots, \mathcal{V}_{n,j}) = 1] - \Pr[D(V_{A,j}, V_{1,j}, \dots, V_{n,j}) = 1]| \\ \leq \text{negl}(\kappa) \end{aligned}$$

as required. We conclude that Protocol 3 securely realizes  $\mathcal{F}_{\text{Time-Delay-MPC}}$  clocked by  $\mathcal{A}'_{i^*}$ .  $\square$

A few remarks are in order. In Protocol 3, all the parties can stop interacting as soon as all the puzzles are outputted. When the locking algorithm  $\text{Lock}(d, t)$  has run-time that is independent of the delay  $t$ , the run-time of Protocol 3 is also independent of the delay parameters. (This is achievable using the [RSW96] time-lock construction, for example.) Alternatively, using a single time-line puzzle in place of the time-lock puzzles in Protocol 3 can improve efficiency, since the time required to generate a time-line puzzle is dependent only on the longest delay  $t_n$ , whereas the time required to generate  $n$  separate time-lock puzzles depends on the sum of all the delays,  $t_1 + \dots + t_n$ .

## 4.4 Time-line puzzles

We now introduce the more general, novel definition of *time-line* puzzles, which can be useful for locking together many data items with different delays for a single recipient, or for locking data for a group of people. In the latter case, it becomes a concern that computation speed will vary between parties: indeed, the scheme will be unworkable if some parties have orders of magnitude more computing power than others, so some assumption is required on the similarity of computing power among parties, such as the

Relative-Delay Assumption of Section 4.3.2. When a time-line puzzle is given to a single recipient, then no additional assumptions are required.

We remark that time-line puzzles could be used (instead of a set of time-lock puzzles) to realize Protocol 3. More generally, we present this new notion because we believe that time-line puzzles may be of independent interest as a timed-release primitive.

In some ways, a time-line puzzle can be thought of as a primitive that packages a sequence of time-lock puzzles together into a unified system about which we can reason and give security guarantees. However, time-line puzzles can also provide concrete advantages over a collection of time-lock puzzles. For example, when issuing many time-lock puzzles to one recipient, the recipient has to run the computation for all of the puzzles in parallel: that is, he does  $O(m \cdot t)$  computation where  $m$  is the number of data items and  $t$  is the time-delay. If instead he gets a time-line puzzle, he only has to run one puzzle's worth of computation in order to unlock all the data items: that is, he does only  $O(t)$  computation, just like for a single time-lock puzzle.

**DEFINITION 16 (TIME-LINE PUZZLES).** A time-line puzzle scheme is a family of PPT algorithms

$$\mathcal{T} = \{(\text{Lock}_m, \text{TimeStep}_m, \text{Unlock}_m)\}_{m \in \mathbb{N}}$$

as follows:

- $\text{Lock}_m(1^\kappa, (d_1, \dots, d_m), (t_1, \dots, t_m))$  takes parameters  $\kappa \in \mathbb{N}$  the security parameter,  $(d_1, \dots, d_m) \in (\{0, 1\}^\ell)^m$  the data items to be locked, and  $(t_1, \dots, t_m) \in \mathbb{N}^m$  the number of steps needed to unlock each data item (respectively), and outputs a puzzle

$$P = (x, (t_1, \dots, t_m), (b_1, \dots, b_m), a)$$

$$\in \{0, 1\}^n \times \mathbb{N} \times (\{0, 1\}^{n'})^m \times \{0, 1\}^{n'}$$

where  $n, n', n'' = \text{poly}(\kappa)$ , and  $a$  can be thought of as auxiliary information.

- $\text{TimeStep}_m(1^\kappa, x', a')$  takes parameters  $\kappa \in \mathbb{N}$  the security parameter, a bit-string  $x' \in \{0, 1\}^n$ , and auxiliary information  $a'$ , and outputs a bit-string  $x'' \in \{0, 1\}^{n'}$ .
- $\text{Unlock}_m(1^\kappa, x', b')$  takes parameters  $\kappa \in \mathbb{N}$  the security parameter, a bit-string  $x' \in \{0, 1\}^n$ , and auxiliary information  $b' \in \{0, 1\}^{n'}$ , and outputs some data  $d' \in \{0, 1\}^\ell$ .

In terms of the “pointer chain” analogy above, solving a time-line puzzle may be thought of as following a pointer chain where not one but many keys are placed along the chain, at different locations  $t_1, \dots, t_m$ . Each key  $x_{t_i}$  in the pointer chain depicted below enables the “unlocking” of the locked data  $b_i$ : for example,  $b_i$  could be the encryption of the  $i^{\text{th}}$  data item  $d_i$  under the key  $x_{t_i}$ .

Using similar notation to that defined for time-lock puzzles: for a time-line puzzle scheme  $\mathcal{T}$ , let  $\text{IterateTimeStep}_m^{\mathcal{T}}$  denote the following function:

$$\text{IterateTimeStep}_m^{\mathcal{T}}(i, x, a) = \underbrace{\text{TimeStep}_m(\text{TimeStep}_m(\dots(\text{TimeStep}_m(x, a), a) \dots), a)}_i.$$

Define  $\text{CompleteUnlock}_{m,i}^{\mathcal{T}}$  to be the following function:

$$\begin{aligned} \text{CompleteUnlock}_{m,i}^{\mathcal{T}}((x, t_i, b_i, a)) = \\ \text{Unlock}_m(\text{IterateTimeStep}_m^{\mathcal{T}}(t_i, x, a), b_i), \end{aligned}$$

that is, the function that should be used to unlock the  $i^{\text{th}}$  piece of data locked by a time-line puzzle which was generated by  $\text{Lock}_m$ . We now define correctness and security for time-line puzzle schemes.

**DEFINITION 17 (CORRECTNESS).** *A time-line puzzle scheme  $\mathcal{T}$  is correct if for all  $m = \text{poly}(\kappa)$  and for all  $i \in [m]$ , it holds that*

$$\Pr_{(x, \vec{t}, \vec{b}, a) \leftarrow \text{Lock}_m(\vec{d}, \vec{t})} [\text{CompleteUnlock}_i^{\mathcal{T}}((x, t_i, b_i, a)) \neq d_i] \leq \text{negl}(\kappa),$$

where  $\kappa$  is the security parameter,  $\vec{d} = (d_1, \dots, d_m)$ , and  $\vec{t} = (t_1, \dots, t_m)$ .

Security for time-line puzzles involves more stringent requirements than security for time-lock puzzles. We define security in terms of two properties which must be satisfied: *timing* and *hiding*. The timing property is very similar to the security requirement for time-lock puzzles, and gives a guarantee about the relative amounts of time required to solve different time-lock puzzles. The hiding property ensures (informally speaking) that the ability to unlock any given data item that is locked in a time-line puzzle does not imply the ability to unlock any others. The security definition (Definition 18, below) refers to the following security experiment.

The experiment  $\text{HidingExp}_{\mathcal{A}, \mathcal{T}}(\kappa)$

1.  $\mathcal{A}$  outputs  $m = \text{poly}(\kappa)$  and data vectors  $\vec{d}_0, \vec{d}_1 \in (\{0, 1\}^\ell)^m$  and a time-delay vector  $\vec{t} \in \mathbb{N}^m$ .
2. The challenger samples  $(\beta_1, \dots, \beta_m) \leftarrow \{0, 1\}^m$ , computes the time-line puzzle  $(x, \vec{t}, \vec{b}, a) = \text{Lock}_m(1^\kappa, ((d_{\beta_1})_1, \dots, (d_{\beta_m})_m), \vec{t})$ , and sends  $(x, a)$  to  $\mathcal{A}$ .
3.  $\mathcal{A}$  sends a query  $i \in [m]$  to the challenger. The challenger responds by sending  $b_i$  to  $\mathcal{A}$ . This step may be repeated up to  $m - 1$  times. Let  $I$  denote the set of queries made by  $\mathcal{A}$ .
4.  $\mathcal{A}$  outputs  $i' \in [m]$  and  $\beta' \in \{0, 1\}$ .
5. The output of the experiment is 1 if  $i' \notin I$  and  $\beta' = \beta_{i'}$ . Otherwise, the output is 0.

**DEFINITION 18 (SECURITY).** *Let*

$$\mathcal{T} = \{(\text{Lock}_m, \text{TimeStep}_m, \text{Unlock}_m)\}_{m \in \mathbb{N}}$$

*be a time-line puzzle scheme.  $\mathcal{T}$  is secure if it satisfies the following two properties.*

- **TIMING:** *For all  $m = \text{poly}(\kappa)$  and  $\vec{d}, \vec{d}' \in (\{0, 1\}^\ell)^m$  and  $\vec{t} = (t_1, \dots, t_m)$ , if there exists an adversary  $\mathcal{A}$  that solves any one of the puzzles defined by the time-line, that is,*

$$\Pr_{P \leftarrow \text{Lock}_m(\vec{d}, \vec{t})} [\mathcal{A}(P) = d_i] = \varepsilon$$

*for some non-negligible  $\varepsilon$  and some  $i \in [m]$ , then for all  $j \in [t_i]$  and all  $\vec{t}' \in [t_m]^m$ , there exists an adversary  $\mathcal{A}_{j, \vec{t}'}$  such that*

$$\Pr_{P' \leftarrow \text{Lock}_m(\vec{d}', \vec{t}')} [\mathcal{A}_{j, \vec{t}'}(P') = d_j] \geq 1 - \text{negl}(\kappa), \text{ and}$$

$$\Pr_{\substack{P \leftarrow \text{Lock}_m(\vec{d}, \vec{t}) \\ P' \leftarrow \text{Lock}_m(\vec{d}', \vec{t}')}} [\text{time}_{\mathcal{A}}(P) \geq (t'_j/t_i) \cdot \text{time}_{\mathcal{A}_{j, \vec{t}'}}(P')]$$

$$| \mathcal{A}(\text{Lock}_m(\vec{d}, \vec{t})) = d_i | \geq 1 - \text{negl}(\kappa).$$

- **HIDING:** *For all PPT adversaries  $\mathcal{A}$ , it holds that*

$$\Pr[\text{HidingExp}_{\mathcal{A}, \mathcal{T}}(\kappa) = 1] \leq 1/2 + \text{negl}(\kappa).$$

In Appendix ??, we describe and prove the security of two constructions of time-line puzzle schemes. One of these schemes is based on a concrete assumption (specifically, on the sequentiality of modular exponentiation, like the time-lock puzzles of [RSW96]), whereas the other is based on the existence of a “black-box” inherently-sequential hash function.

## Acknowledgements

We would like to thank Yehuda Lindell for an interesting discussion on the nature of fairness in multiparty computation, and we are grateful to Juan Garay, Björn Tackmann, and Vassilis Zikas for an illuminating discussion about measures of partial fairness in MPC.

The first author acknowledges research support by the Robert Solow Fellowship 3310100. The second and third authors acknowledge research supported by NSF Eager CNS-1347364, NSF Frontier CNS-1413920, the Simons Foundation (agreement dated June 5, 2012), Air Force Laboratory FA8750-11-2-0225, and Lincoln Lab PO7000261954. This work was done in part while these authors were visiting the Simons Institute for the Theory of Computing, supported by the Simons Foundation and by the DIMACS/Simons Collaboration in Cryptography through NSF grant CNS-1523467.

## References

- [AM12] Pablo Daniel Azar and Silvio Micali. “Rational proofs”. In: *Proceedings of the forty-fourth annual ACM symposium on Theory of computing*. ACM. 2012, pp. 1017–1028.
- [AM13] Pablo Daniel Azar and Silvio Micali. “Super-efficient rational proofs”. In: *Proceedings of the fourteenth ACM conference on Electronic commerce*. ACM. 2013, pp. 29–30.
- [BGK14] Siddhartha Banerjee, Ashish Goel, and Anilesh Kollagunta Krishnaswamy. “Re-incentivizing discovery: mechanisms for partial-progress sharing in research”. In: *ACM Conference on Economics and Computation, EC '14, Stanford, CA, USA, June 8-12, 2014*. Ed. by Moshe Babaioff, Vincent Conitzer, and David Easley. ACM, 2014, pp. 149–166. ISBN: 978-1-4503-2565-3. DOI: 10.1145/2600057.2602888. URL: <http://doi.acm.org/10.1145/2600057.2602888>.

- [BL+07] Yoshua Bengio, Yann LeCun, et al. "Scaling learning algorithms towards AI". In: *Large-scale kernel machines* 34.5 (2007).
- [BK14] Iddo Bentov and Ranjit Kumaresan. "How to Use Bitcoin to Design Fair Protocols". In: *Advances in Cryptology - CRYPTO 2014 - 34th Annual Cryptology Conference, Santa Barbara, CA, USA, August 17-21, 2014, Proceedings, Part II*. Ed. by Juan A. Garay and Rosario Gennaro. Vol. 8617. Lecture Notes in Computer Science. Springer, 2014, pp. 421–439. ISBN: 978-3-662-44380-4. DOI: 10.1007/978-3-662-44381-1\_24. URL: [http://dx.doi.org/10.1007/978-3-662-44381-1\\_24](http://dx.doi.org/10.1007/978-3-662-44381-1_24).
- [BP12] Sarah E. Bergen and Tracey L. Petryshen. "Genome-wide association studies (GWAS) of schizophrenia: does bigger lead to better results?" In: *Current opinion in psychiatry* 25.2 (Mar. 2012), pp. 76–82.
- [BN00] Dan Boneh and Moni Naor. "Timed Commitments". In: *Advances in Cryptology - CRYPTO 2000, 20th Annual International Cryptology Conference, Santa Barbara, California, USA, August 20-24, 2000, Proceedings*. Ed. by Mihir Bellare. Vol. 1880. Lecture Notes in Computer Science. Springer, 2000, pp. 236–254. ISBN: 3-540-67907-3. DOI: 10.1007/3-540-44598-6\_15. URL: [http://dx.doi.org/10.1007/3-540-44598-6\\_15](http://dx.doi.org/10.1007/3-540-44598-6_15).
- [CDP14] Yang Cai, Constantinos Daskalakis, and Christos Papadimitriou. "Optimum Statistical Estimation with Strategic Data Sources". In: *ArXiv e-prints* (2014). arXiv: 1408.2539 [stat.ML].
- [Can00] Ran Canetti. "Security and Composition of Multiparty Cryptographic Protocols". In: *J. Cryptology* 13.1 (2000), pp. 143–202. DOI: 10.1007/s001459910006. URL: <http://dx.doi.org/10.1007/s001459910006>.
- [Cle86] Richard Cleve. "Limits on the Security of Coin Flips when Half the Processors Are Faulty (Extended Abstract)". In: *Proceedings of the 18th Annual ACM Symposium on Theory of Computing, May 28-30, 1986, Berkeley, California, USA*. Ed. by Juris Hartmanis. ACM, 1986, pp. 364–369. ISBN: 0-89791-193-8. DOI: 10.1145/12130.12168. URL: <http://doi.acm.org/10.1145/12130.12168>.
- [Edm65] Jack Edmonds. "Paths, trees, and flowers". In: *Canadian Journal of mathematics* 17.3 (1965), pp. 449–467.
- [ENSS95] Guy Even, Joseph (Seffi) Naor, Baruch Schieber, and Madhu Sudan. "Approximating minimum feedback sets and multi-cuts in directed graphs". English. In: *Integer Programming and Combinatorial Optimization*. Ed. by Egon Balas and Jens Clausen. Vol. 920. Lecture Notes in Computer Science. Springer Berlin Heidelberg, 1995, pp. 14–28. ISBN: 978-3-540-59408-6. DOI: 10.1007/3-540-59408-6\_38. URL: [http://dx.doi.org/10.1007/3-540-59408-6\\_38](http://dx.doi.org/10.1007/3-540-59408-6_38).
- [For11] Schizophrenia Research Forum. *GWAS Goes Bigger: Large Sample Sizes Uncover New Risk Loci, Additional Overlap in Schizophrenia and Bipolar Disorder*. Sept. 2011. URL: <http://www.schizophreniaforum.org/new/detail.asp?id=1692>.
- [GMW87] Oded Goldreich, Silvio Micali, and Avi Wigderson. "How to Play any Mental Game or A Completeness Theorem for Protocols with Honest Majority". In: *Proceedings of the 19th Annual ACM Symposium on Theory of Computing, 1987, New York, New York, USA*. Ed. by Alfred V. Aho. ACM, 1987, pp. 218–229. ISBN: 0-89791-221-7. DOI: 10.1145/28395.28420. URL: <http://doi.acm.org/10.1145/28395.28420>.
- [GHRV14] Siyao Guo, Pavel Hubáček, Alon Rosen, and Margarita Vald. "Rational arguments: single round delegation with sublinear verification". In: *Proceedings of the 5th conference on Innovations in theoretical computer science*. ACM, 2014, pp. 523–540.
- [Han12] Robin Hanson. "Logarithmic Market Scoring Rules For Modular Combinatorial Information Aggregation". In: *The Journal of Prediction Markets* 1.1 (2012), pp. 3–15.
- [Ins14] Broad Institute. *International team sheds new light on biology underlying schizophrenia*. July 2014. URL: <https://www.broadinstitute.org/news/5895>.
- [Kar72] Richard M Karp. *Reducibility among combinatorial problems*. Springer, 1972.
- [KO11] Jon M. Kleinberg and Sigal Oren. "Mechanisms for (mis)allocating scientific credit". In: *Proceedings of the 43rd ACM Symposium on Theory of Computing, STOC 2011, San Jose, CA, USA, 6-8 June 2011*. Ed. by Lance Fortnow and Salil P. Vadhan. ACM, 2011, pp. 529–538. ISBN: 978-1-4503-0691-1. DOI: 10.1145/1993636.1993707. URL: <http://doi.acm.org/10.1145/1993636.1993707>.
- [May93] Timothy C. May. *Timed-release crypto*. 1993. URL: <http://www.hks.net/cpunks/cpunks-01460.html>.
- [PGC14] Schizophrenia Working Group of the Psychiatric Genomics Consortium. "Biological insights from 108 schizophrenia-associated genetic loci". In: *Nature* 511 (Mar. 2014), pp. 421–427.
- [RSW96] Ronald L. Rivest, Adi Shamir, and David A. Wagner. *Time-lock puzzles and timed-release crypto*. Tech. rep. 1996.
- [Sha79] Adi Shamir. "How to Share a Secret". In: *Commun. ACM* 22.11 (Nov. 1979), pp. 612–613. ISSN: 0001-0782.
- [Wel03] The Wellcome Trust. *Sharing Data from Large-scale Biological Research Projects: A System of Tripartite Responsibility*. Report of a meeting organized by the Wellcome Trust and held on 14–15 January 2003 at Fort Lauderdale, USA. 2003.