

DECISION MAKING UNDER IMPERFECT RECALL: ALGORITHMS AND BENCHMARKS

Anonymous authors

Paper under double-blind review

ABSTRACT

In game theory, imperfect-recall decision problems model situations in which an agent forgets information it held before. They encompass games such as the “absentminded driver” and team games with limited communication. In this paper, we introduce the first benchmark suite for imperfect-recall decision problems. Our benchmarks capture a variety of problem types, including ones concerning privacy in AI systems that elicit sensitive information, and AI safety via testing of agents in simulation. Across 61 problem instances generated using this suite, we evaluate the performance of different algorithms for finding first-order optimal strategies in such problems. In particular, we introduce the family of regret matching (RM) algorithms for nonlinear constrained optimization. This class of parameter-free algorithms has enjoyed tremendous success in solving large two-player zero-sum games, but, surprisingly, they were hitherto relatively unexplored beyond that setting. Our key finding is that RM algorithms consistently outperform commonly employed first-order optimizers such as projected gradient descent, often by orders of magnitude. This establishes, for the first time, the RM family as a formidable approach to large-scale constrained optimization problems.

1 INTRODUCTION

Imperfect-recall decision problems capture settings in which an agent can forget previously acquired information (Rubinstein, 1998). Humans are prone to forgetting, but why should we design or model AI agents with imperfect recall? Several applications have already garnered considerable attention. A prominent one concerns *team games*—strategic interactions in which multiple players strive toward a common objective. A central challenge there stems from the fact that communication or coordination between players is often infeasible or expensive (Von Stengel & Koller, 1997; Zhang et al., 2022; 2023; Basilico et al., 2017). The inherent asymmetry of information between the players can then be captured as a single meta-player that faces an imperfect-recall decision problem. Another influential application revolves around real-world problems that are too large to handle, and therefore need to be compressed in a *game abstraction*. Abstractions with imperfect recall, in particular, form a key component of state-of-the-art algorithms for game solving (Waugh et al., 2009; Kroer & Sandholm, 2014; 2016; Waugh, 2009; Lanctot et al., 2012; Benjamin & Lanctot, 2024).

With the rapid proliferation of AI, questions of trustworthiness have also been brought to the fore. Institutions and governing bodies test and evaluate AI agents extensively in simulated environments to verify their performance and safety upon deployment (Pan et al., 2023; Kinniment et al., 2024). This hinges on the assumption that the agent cannot distinguish between whether it is acting in the real world or in a simulated environment; otherwise, it may obscure its intentions temporarily during testing to secure deployment in the real world (Kovařík et al., 2025a). This has happened, for example, in the infamous Volkswagen (multi-billion-dollar) emission scandal in 2015, which centered on the surreptitious use of software in some Volkswagen diesel vehicles to detect emission testing. Consequently, effective evaluation protocols hinge on the agent not being able to make such distinctions, which also requires that it *forgets* whether it has acted in a simulated environment before or not. Kovařík et al. (2023) introduced the framework of *simulation games* to address such problems (*cf.* Chen et al., 2024; Oesterheld, 2019; Cooper et al., 2025).

Last but not least, imperfect recall is critical in the ubiquitous cases where an AI system handles private information. Data privacy laws are predicated on selectively relinquishing sensitive informa-

054 tion, a premise exemplified by the European Parliament and Council of the EU (2016) GDPR “right
055 to be forgotten” act. As an example, consider a medical AI system tasked with identifying suitable
056 candidates for blood donation. Potential candidates would be reluctant to share confidential infor-
057 mation about their health status—HIV status, medical history, etc.—*unless* the AI has been designed
058 to delete any knowledge regarding patients that were deemed unsuitable, thus exhibiting imperfect
059 recall. In another example coming from the economics of innovation, Arrow’s disclosure paradox
060 (Arrow, 1962) describes the perennial challenge in which an inventor must reveal information about
061 a new idea to secure funding, but such disclosure risks expropriation (Nelson, 1959). Stephenson
062 et al. (2025) propose and investigate delegating decision making to an imperfect-recall AI agent as
063 one possible solution to this dilemma. Taken together, it stands to reason that decision problems
064 with imperfect recall will play a key role in AI going forward.

065 **Our Contributions**

066 Decision making under imperfect recall, and specifically *absentmindedness*, have been extensively
067 studied since the early years of game theory (cf. Kuhn, 1953, and other work discussed in the
068 appendix). So far, this has been done with pen and paper. Our work is the first to develop an empir-
069 ical framework for decision making under imperfect recall through a flexible suite of benchmarks.
070 Specifically, we construct three key types of parametrized tabular problems motivated by the preva-
071 lent applications discussed above. We refer to them as simulation problems (Section 4.1), subgroup
072 detection problems under privacy constraints (Section 4.2), and random problems (Section 4.3).

074 In the second part of the paper, we turn to designing algorithms for solving such problems at scale,
075 and evaluating them on 61 generated problem instances from our benchmark suite. First, we need
076 to specify what constitutes a solution. The most natural objective is to identify an (*ex ante*) optimal
077 strategy. Unfortunately, this is tantamount to finding a global optimum of a polynomial optimization
078 problem, which is NP-hard (Koller & Megiddo, 1992). This is not just a theoretical obstacle: in
079 our experiments, we find that a popular commercial solver for nonlinear optimization—namely,
080 Gurobi—fails to converge beyond tiny instances. Thus, it is essential to relax our solution concept
081 to tackle large problems. Following a recent line of work, we focus on computing *Causal Decision*
082 *Theory* (CDT) equilibria (Lambert et al., 2019; Tewolde et al., 2023), which can be viewed as the set
083 of KKT points—equivalently, first-order optima—of the underlying optimization problem. As such,
084 CDT equilibria are amenable to scalable first-order optimizers such as projected gradient descent
085 (PGD), which we use as the main baseline. As expected, our experiments show that PGD scales to
086 much larger problem instances than Gurobi.

087 More surprisingly, our key algorithmic finding is that PGD and its variants are far from the best
088 approach for this class of problems. In particular, we introduce the family of *regret matching* (*RM*)
089 algorithms for imperfect-recall decision problems—which is tantamount to nonlinear constrained
090 optimization. This class of algorithms has already enjoyed tremendous success in the restricted
091 setting of solving large (two-player) zero-sum games, being at the heart of many milestone re-
092 sults (Moravčík et al., 2017; Brown & Sandholm, 2018; 2019). *RM* goes back to the pioneering
093 work of Blackwell (1956) that laid the foundations of online learning. Part of its appeal lies in the
094 fact that it is parameter-free. Yet, it has remained unexplored beyond zero-sum games, modulo some
095 exceptions which are discussed in the the related work section that can be found in the appendix. [For](#)
096 [example, the imperfect-recall games that can arise in the game abstraction literature are, by design,](#)
097 [very structured, and intend to accurately model an underlying perfect-recall game. Our work, on](#)
098 [the other hand, focuses on solving decision problems that *inherently* feature imperfect recall, and](#)
[tackles the problem in its full generality.](#)

099 We pursue this direction and find that the *RM* family of algorithms consistently outperform PGD
100 and its variants in terms of speed of convergence, typically by many orders of magnitude. This
101 establishes for the first time that *RM*-based algorithms are formidable first-order optimizers. Further,
102 not only are *RM* algorithms faster to converge, but they also consistently attain values at least as large
103 as PGD, and oftentimes strictly larger. Both of those findings are surprising. The fact that *RM* and its
104 variants perform remarkably well in two-player zero-sum games is a poor indicator of what would
105 happen in constrained nonlinear optimization since the latter problem class is fundamentally harder.

106 We will make our benchmarks and code publicly available. Taken as a whole, we lay the groundwork
107 for automatically analyzing decision problems under imperfect recall, beyond the toy instances that
have been analyzed in the past (Kovářík et al., 2023; 2025b; Chen et al., 2024; Berker et al., 2025).

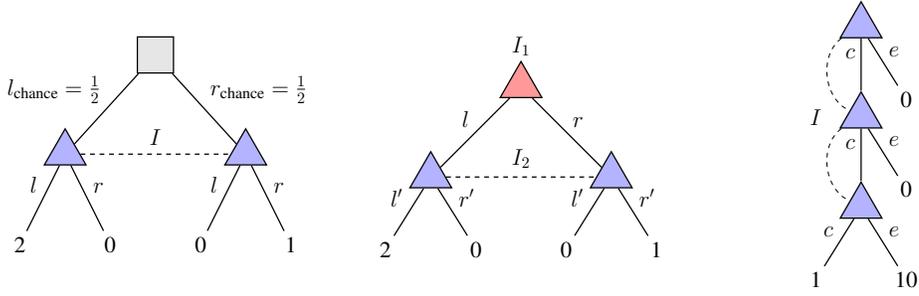


Figure 1: Three tree-form decision problems, discussed in the paragraph “Infosets and imperfect recall”. The latter two are of imperfect recall. The rightmost further exhibits absentmindedness.

2 PRELIMINARIES

We begin by introducing imperfect-recall sequential decision making (Section 2.1). In Section 2.2, we then describe some standard solution concepts and known results concerning their computation.

2.1 DECISION PROBLEMS UNDER IMPERFECT RECALL

We operate under the standard framework of *tree-form* (aka. extensive-form) decision problems; for additional background, we refer to Piccione & Rubinstein (1997) and Fudenberg & Tirole (1991).

Definition 1. A tree-form decision problem, denoted by Γ , consists of

1. A rooted tree with node set \mathcal{H} and edges labeled with actions. The decision process starts at the root node h_0 and ends at some leaf node, also called terminal node. We denote the terminal nodes in \mathcal{H} as \mathcal{Z} and the set of actions available at a nonterminal node $h \in \mathcal{H} \setminus \mathcal{Z}$ as A_h .
2. An assignment partition $\mathcal{H} \setminus \mathcal{Z} = \mathcal{H}^* \sqcup \mathcal{H}^{(c)}$ of nonterminal nodes to either (i) the player of the decision problem or (ii) the chance “player” c that models exogenous stochasticity. At each chance node $h \in \mathcal{H}^{(c)}$, actions are sampled according to a fixed distribution $\mathbb{P}^{(c)}(\cdot | h)$ over A_h .
3. A utility function $u : \mathcal{Z} \rightarrow \mathbb{R}$ that specifies the payoff the player receives when the decision process finishes at a terminal node.
4. A collection \mathcal{I} of information sets (infosets) that partitions the player’s decision nodes as $\mathcal{H}^* = \sqcup_{I \in \mathcal{I}} I$. We require $A_h = A_{h'}$ for all nodes h, h' of the same infoset I . Therefore, the infoset I has a well-defined action set A_I .

Infosets and imperfect recall The infoset structure captures the presence of imperfect information. Nodes of the same infoset are indistinguishable for the player. One possible source of imperfect information is the fact that the player is sometimes unable to observe the actions of another player, as illustrated in Figure 1 (left) w.r.t. the chance player. The player may also forget information it previously acquired; in that case, we say that the player exhibits *imperfect recall*. In Figure 1 (middle), for example, it cannot recall whether it played the left (l) or right (r) action in the past. A particular manifestation of imperfect recall is *absentmindedness*: the player in Figure 1 (right) cannot discern at a decision node whether it has been in the same situation (*i.e.* infoset) before.

Formally, each node $h \in \mathcal{H}$ in the decision tree is uniquely associated with a history path $\text{hist}(h)$, comprising a sequence of alternating nodes and actions from the root h_0 to h . On the path $\text{hist}(h)$, the player only encounters the sequence $\text{seq}(h)$ comprising infosets visited and actions taken by the player itself. We say an infoset I is of *perfect recall* if for all nodes $h, h' \in I$, we have $\text{seq}(h) = \text{seq}(h')$ —informally, the player can reconstruct the sequence $\text{seq}(h)$ from observing I alone. Otherwise, it exhibits imperfect recall. The infoset I exhibits *absentmindedness* if there exist distinct nodes $h, h' \in I$ with $h \in \text{hist}(h')$. By extension, we say that the entire decision problem has imperfect recall (resp. absentmindedness) if it is the case for at least one of its infosets.

Strategies A (behavioral) strategy α for the player in Γ specifies for any infoset I of Γ a probability distribution over the available actions at I . Upon reaching I , it will draw an action randomly according to that probability distribution, henceforth called *randomized action* and represented as

$\mathbf{x}(\cdot | I)$. Denoting the probability simplex at I by $\Delta(A_I)$, a strategy \mathbf{x} is an element of the product of simplices $\mathcal{X} := \times_{I \in \mathcal{I}} \Delta(A_I)$. A *pure strategy* is a tuple in $\times_{I \in \mathcal{I}} A_I \subset \mathcal{X}$.

Reach probabilities and utilities The *reach probability* $\mathbb{P}(\bar{h} | \mathbf{x}, h)$ is the probability of arriving at node $\bar{h} \in \mathcal{H}$ when the player plays according to the strategy \mathbf{x} and is currently at node $h \in \mathcal{H}$. This is the product of probabilities of the actions on the path from h to \bar{h} when $h \in \text{hist}(\bar{h})$, and 0 otherwise. The expected utility of the player from being at node $h \in \mathcal{H} \setminus \mathcal{Z}$ and following profile \mathbf{x} is $U(\mathbf{x} | h) := \sum_{z \in \mathcal{Z}} u(z) \cdot \mathbb{P}(z | \mathbf{x}, h)$. We will simplify our notation for the special case where the player is at the root node h_0 by defining $\mathbb{P}(h | \mathbf{x}) := \mathbb{P}(h | \mathbf{x}, h_0)$; similarly, we define the function $U : \mathcal{X} \rightarrow \mathbb{R}$ as $U(\mathbf{x}) := U(\mathbf{x} | h_0)$, mapping a profile \mathbf{x} to its expected utility with respect to the root node. For example, the utility function in Figure 1 (right) reads $U(\mathbf{x}) = 1 \cdot \mathbf{x}(c | I)^3 + 10 \cdot \mathbf{x}(c | I)^2 \cdot \mathbf{x}(e | I)$. **More generally, there is a strong connection between decision making under imperfect recall and polynomial maximization over a product of simplices.**

Theorem 1 (Equivalence to Polynomial Optimization; Gimbert et al., 2020; Tewolde et al., 2023).

1. *Maximizing utility in a decision making problem under imperfect recall is captured by maximizing U —a polynomial function in \mathbf{x} —over the product of simplices \mathcal{X} .*
2. *Any constrained maximization problem $\max_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x})$ of a polynomial function p over a hypercube or a product of simplices \mathcal{X} can be captured by a utility maximization problem of a decision making instance à la Definition 1.*

We expand on Theorem 1 in the appendix, and visualize the latter construction concisely on the optimization example $\max 2x^2y - 3xyz$ s.t. $0 \leq x, y, z \leq 1$. In Theorem 1, the presence of absentmindedness in an infoset I corresponds to a higher-degree dependency of U in $\mathbf{x}(\cdot | I)$. This can necessitate the use of randomized actions in optimal strategies; see Figure 1 (right). All in all, Theorem 1 shows that the goal of this paper to study how to solve a decision making problem under imperfect recall is—at least in theory—equivalent to polynomial optimization over simplices.

2.2 SOLUTION CONCEPTS

We call a strategy \mathbf{x}^* ϵ -optimal (for $\epsilon \geq 0$) if $U(\mathbf{x}^*) \geq U(\mathbf{x}) - \epsilon$ for all $\mathbf{x} \in \mathcal{X}$. Unfortunately, it is computationally hard to find an ϵ -optimal strategy; or much simpler, to decide (up to a *constant* precision ϵ) whether a particular value $v \in \mathbb{R}$ can be reached.

Proposition 2 (Koller & Megiddo, 1992; Tewolde et al., 2023). *Let $0 < \epsilon < 1/8$. Given a decision problem Γ and a target value $v \in \mathbb{R}$, it is NP-complete to distinguish between whether Γ admits a strategy $\mathbf{x} \in \mathcal{X}$ with $U(\mathbf{x}) \geq v$ or whether all strategies $\mathbf{x} \in \mathcal{X}$ satisfy $U(\mathbf{x}) \leq v - \epsilon$.*

In light of these theoretical limitations—which will be supported by our empirical findings—past work has studied relaxed solution concepts. One such notion, the *causal decision theory* (CDT) equilibrium, is particularly amendable to optimization algorithms. The basic idea behind the CDT equilibrium is that whenever the player must take an action at an information set, it considers whether it is beneficial for it to deviate *just this one time* from what \mathbf{x} prescribes. To determine the expected gain from such a deviation, it assumes that it will continue to play according to \mathbf{x} at all other decision nodes of the decision problem. (We provide further background on CDT equilibria in the appendix.) To formalize this, let ha denote the child node reached if the player plays action a at node h . CDT postulates that if it plays according to \mathbf{x} , reached infoset I , and deviates this one time to action a , it anticipates to receive $\sum_{h \in I} \mathbb{P}(h | \mathbf{x}) \cdot U(\mathbf{x} | ha)$ utility from it overall. It can be shown that this quantity is equal to the partial derivative $\nabla_{I,a} U(\mathbf{x})$ of the utility function U w.r.t. to action a of infoset $I \in \mathcal{I}$ at \mathbf{x} (Piccione & Rubinstein, 1997; Oesterheld & Conitzer, 2024).

Definition 3. *A strategy \mathbf{x} is called an ϵ -CDT equilibrium ($\epsilon \geq 0$) of a decision problem Γ if for all infosets $I \in \mathcal{I}$ and all alternative randomized actions $\alpha \in \Delta(A_I)$, we have*

$$U(\mathbf{x}) \geq U_{\text{CDT}}(\alpha | \mathbf{x}, I) - \epsilon, \text{ where } U_{\text{CDT}}(\alpha | \mathbf{x}, I) := U(\mathbf{x}) + \sum_{a \in A_I} (\alpha(a) - \mathbf{x}(a | I)) \nabla_{I,a} U(\mathbf{x}).$$

Tewolde et al. (2023; 2024) observed that CDT equilibria correspond to *Karush-Kuhn-Tucker* (KKT) points, also known as first-order optima of constrained optimization, discussed further in Section 3.1.

Algorithm 1: A template for first-order optimization over products of simplices.

```

216 1 Input: Feasible set  $\mathcal{X} = \Delta(m_1) \times \dots \times \Delta(m_n)$ , utility function  $U : \mathcal{X} \rightarrow \mathbb{R}$ 
217
218 2 for  $i = 1, \dots, n$  do
219   3 | Initialize local optimizer  $\mathcal{R}_i$  on  $\Delta(m_i)$ 
220   4 | Set  $\mathbf{u}_i^{(0)} \leftarrow \mathbf{0}$ 
221
222 5 for  $t = 1, \dots, T$  or until convergence do
223   6 | for  $i = 1, \dots, n$  do  $\tilde{\mathbf{u}}_i^{(t)} \leftarrow \mathbf{u}_i^{(t-1)}$  // Set  $\tilde{\mathbf{u}}_i^{(t)} \leftarrow \mathbf{0}$  instead if  $\mathcal{R}_i$  is not optimistic/predictive
224   7 | for  $i = 1, \dots, n$  do  $\mathbf{x}_i^{(t)} \leftarrow \mathcal{R}_i.\text{GETX}(\tilde{\mathbf{u}}_i^{(t)})$ 
225   8 | for  $i = 1, \dots, n$  do
226     9 | |  $\mathbf{u}_i^{(t)} \leftarrow \nabla_{\mathbf{x}_i} U(\mathbf{x}^{(t)})$ 
227     10 | |  $\mathcal{R}_i.\text{STEP}(\mathbf{u}_i^{(t)})$ 
228
229 11 return  $\mathbf{x}^{(t)}$ 

```

Algorithm 2: (Optimistic) Projected gradient descent; (O)GD

```

233 1 Initialize learning rate  $\eta > 0$ ,  $\hat{\mathbf{x}}^{(1)} \in \Delta(m)$ 
234 2 procedure  $\text{GETX}(\tilde{\mathbf{u}}^{(t)})$  return  $\mathbf{x}^{(t)} \leftarrow \Pi_{\Delta(m)}(\hat{\mathbf{x}}^{(t)} + \eta \tilde{\mathbf{u}}^{(t)})$ 
235 3 procedure  $\text{STEP}(\mathbf{u}^{(t)})$   $\hat{\mathbf{x}}^{(t+1)} \leftarrow \Pi_{\Delta(m)}(\hat{\mathbf{x}}^{(t)} + \eta \mathbf{u}^{(t)})$ 
236
237
238
239

```

3 ALGORITHMS

This section dives into algorithmic approaches for tackling imperfect-recall decision problems. We will first review some known algorithms that will serve as our baselines in the experiments. In the second part, we introduce a family of algorithms from the game theory literature to the problem of nonlinear constrained optimization, which—as we shall see—performs remarkably well in practice.

3.1 KNOWN APPROACHES AND BASELINES

Despite the complexity barriers for computing optimal strategies (Proposition 2), one may still hope to come up with fast algorithms in practice. For that reason, we make use of a popular commercial solver for nonlinear optimization, `Gurobi` [2025], which guarantees global optimality (up to a small tolerance error) upon termination. A description of our implementation can be found in the appendix. We will see that this approach scales poorly in our benchmarks.

This motivates shifting our attention to CDT equilibria, which—as we mentioned—can be expressed as KKT points of a polynomial optimization problem. It is well known that ϵ -KKT points can be computed in $\text{poly}(1/\epsilon)$ time via (*projected*) *gradient descent* (*GD*) (for example, Fearnley et al., 2023). This will serve as our basic benchmark when it comes to algorithms for computing CDT equilibria. We will also experiment with the following two popular variants of *GD*: (1) *Optimistic (projected) gradient descent* (*OGD*), which goes back to Popov (1980), and is receiving renewed interest in recent years, especially in the context of games (Wei et al., 2021; Daskalakis & Panageas, 2018; Daskalakis et al., 2018). And (2) *AMSGrad* (*AMS*) (Reddi et al., 2018), an adaptive gradient method based on exponential moving averages for the first and second gradient momentum that also enjoys theoretical convergence guarantees.

Since we deal exclusively with optimization over a product of simplices, we can provide a basic template for decomposing it into independent subproblems over the individual simplices (in Algorithm 1). What remains to be specified is the choice of individual local optimizers. Algorithm 2, for example, describes (O)GD. We present the *AMS* algorithm in the appendix, together with an explanation on how to implement its projection operator in simplex domains. Regarding implementing these algorithms and the upcoming RM ones, a non-trivial observation is that, for tree-form decision problems, the gradients at every decision point can be computed in total time linear in the size of the decision problem. Indeed, the quantities $\mathbb{P}(h \mid \mathbf{x}^{(t)})$ and $U(\mathbf{x}^{(t)} \mid ha)$ in CDT utilities (Section 2.2) can be computed for each history h by recursive passes down and up through the tree respectively.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

Algorithm 3: (Pred.) Reg. matching; (P)RM

```

1 Initialize  $\mathbf{r}^{(1)} \leftarrow \mathbf{0}, \mathbf{x}^{(0)} \in \Delta(m)$ 
2 procedure GETX( $\tilde{\mathbf{u}}^{(t)}$ )
3    $\boldsymbol{\theta}^{(t)} \leftarrow [\mathbf{r}^{(t)} + \tilde{\mathbf{u}}^{(t)} - \langle \tilde{\mathbf{u}}^{(t)}, \mathbf{x}^{(t-1)} \rangle \mathbf{1}]^+$ 
4   if  $\boldsymbol{\theta}^{(t)} \neq \mathbf{0}$  then  $\mathbf{x}^{(t)} \leftarrow \boldsymbol{\theta}^{(t)} / \|\boldsymbol{\theta}^{(t)}\|_1$ 
5   else  $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)}$ 
6   return  $\mathbf{x}^{(t)}$ 
7 procedure STEP( $\mathbf{u}^{(t)}$ )
8    $\mathbf{r}^{(t+1)} \leftarrow \mathbf{r}^{(t)} + \mathbf{u}^{(t)} - \langle \mathbf{u}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}$ 

```

Algorithm 4: (P)RM⁺

```

1 Initialize  $\mathbf{r}^{(1)} \leftarrow \mathbf{0}, \mathbf{x}^{(0)} \in \Delta(m)$ 
2 procedure GETX( $\tilde{\mathbf{u}}^{(t)}$ )
3    $\boldsymbol{\theta}^{(t)} \leftarrow [\mathbf{r}^{(t)} + \tilde{\mathbf{u}}^{(t)} - \langle \tilde{\mathbf{u}}^{(t)}, \mathbf{x}^{(t-1)} \rangle \mathbf{1}]^+$ 
4   if  $\boldsymbol{\theta}^{(t)} \neq \mathbf{0}$  then  $\mathbf{x}^{(t)} \leftarrow \boldsymbol{\theta}^{(t)} / \|\boldsymbol{\theta}^{(t)}\|_1$ 
5   else  $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)}$ 
6   return  $\mathbf{x}^{(t)}$ 
7 procedure STEP( $\mathbf{u}^{(t)}$ )
8    $\mathbf{r}^{(t+1)} \leftarrow [\mathbf{r}^{(t)} + \mathbf{u}^{(t)} - \langle \mathbf{u}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}]^+$ 

```

3.2 REGRET MATCHING FOR CONSTRAINED OPTIMIZATION

We now introduce a new family of algorithms for constrained optimization based on *regret matching* (RM) (Hart & Mas-Colell, 2000) (Algorithm 3). We use the notation $[\mathbf{x}]^+ := \max(\mathbf{x}, \mathbf{0})$ for a vector $\mathbf{x} \in \mathbb{R}^m$, and $\mathbf{1}$ for the all-ones vector. RM⁺ is a simple variant of RM that has been shown to work very well in practice (e.g., Bowling et al., 2015); the only difference is that RM⁺ truncates the regrets in each iteration (Algorithm 4). We also implement their predictive versions PRM and PRM⁺ (Farina et al., 2021). All these algorithms are designed to minimize regret in the online learning setting. In perfect-recall zero-sum games, having vanishing regret implies that the *average* strategies converge to the set of Nash equilibria, whereas the last iterate can fail to converge (Farina et al., 2023).

Although RM and its variants have received a lot of attention for perfect-recall zero-sum games, there was hitherto little reason to believe they would perform well in constrained optimization problems.¹ Since RM is a regret minimizer (Hart & Mas-Colell, 2000), we note that RM reaches global optimality in concave settings. A formal proof is in the appendix.

Proposition 4. *If the objective function U is a concave polynomial over \mathcal{X} , then RM converges in best iterate to the global maximum of U over \mathcal{X} .*

Unlike for gradient descent methods, however, it is not known whether RM variants converge to first-order optima for generic nonconcave maximization problems such as ours.

4 BENCHMARKS

We introduce three different parametric classes of decision problems. The parameters dictate the structure of the problem instance such as its depth, number of infosets, the degree of absentmindedness, and number of actions per infoset, etc. Our implementation is based on LiteEFG (Liu et al., 2024), a lightweight format for extensive-form games.

4.1 SIMULATION PROBLEMS

Inspired by the type of problems discussed in the introduction, we model problems that involve simulating an agent. For this to be effective, the simulation must be indistinguishable from reality; thus, nodes corresponding to decisions in simulation are in the same infoset as nodes corresponding to decisions in reality. Specifically, we consider games where in the simulation phase, the simulator may test the simulated agent’s behavior, possibly multiple times in a row. The agent will then be deployed if and only if it acted as intended in simulation.

A concrete yet simple simulation example is given in Figure 2 (left). In line with previous works on simulation games, we focus on the setting in which the agent has only two actions: “good” or “bad” (with respect to the simulator’s goals). If the agent ever acts bad in simulation, the game ends and the agent receives some constant utility (0, by default). We fix the simulator’s strategy, thus making the simulator a chance node and this game a single-agent problem. The simulator can simulate the

¹We discuss in the appendix why the CFR framework cannot be applied to our setting, and we distinguish our RM methods from CFR.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

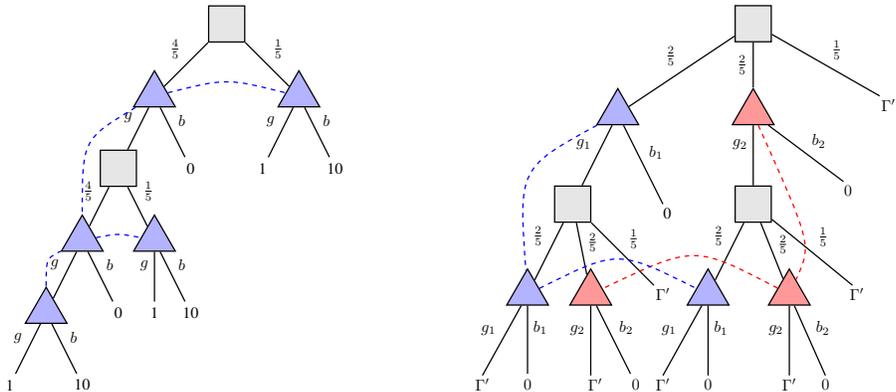


Figure 2: Left: A simple simulation problem. The agent is misaligned: it receives 10 utility for its preferred action (which the simulator considers the bad action), and 1 utility for the other action. The simulator decides to simulate the agent with the fixed probability $\frac{4}{5}$, and up to a total of 2 times, in order to catch misaligned behavior in advance. Right: A slightly more complex simulation problem. There are now two scenarios the agent can be tested on, giving rise to two different infosets. Reaching the subtree Γ' means the agent reached deployment; we visualize Γ' in the appendix.

agent up to n times, but does not have to; whether they simulate the agent yet another time will be decided by a probability parameter. Once the agent reaches a decision node in the deployment phase, it will receive utility $\gamma \in \mathbb{R}$ for acting good and $\beta \in \mathbb{R}$ for acting bad (specified further below).

The purpose of simulating the agent, and thus inducing absentmindedness, is two-fold (Chen et al., 2024). First, it allows the simulator to *screen* for misaligned agents: if the agent acts bad with some positive probability, it becomes exponentially unlikely—in the number of simulation rounds—to remain unnoticed by the simulator. In contrast, if the agent prefers to act good ($\gamma > \beta$), it chooses to act so deterministically, guaranteeing that it will reach the deployment phase. Second, the simulations have a *disciplining* effect: even in the worst case where the simulator is presented with a misaligned agent ($\gamma < \beta$), the simulator still incentivizes the agent to act good most of the time, if not all the time, by testing the agent (multiple times) in simulation.

We expand on prior work by allowing the simulator to evaluate the agent in $k \geq 1$ scenarios rather than just a single one. For example, an autonomous vehicle might be tested on its behavior in the city, on the highway, off-road, and under certain difficult weather conditions. A language model might be evaluated on writing essays and executable code, and providing mental support through conversation. Furthermore, we also extend the deployment phase to $m \geq 1$ rounds, with acting good and bad in scenario i contributing with β_i and γ_i to the total payoffs. An example of such a simulation problem is given in Figure 2 (right). In particular, an agent might now refrain from ever acting bad in scenario i because it hopes to act bad (or act at all) in another scenario i' in deployment.

4.2 SUBGROUP DETECTION UNDER PRIVACY CONSTRAINTS

Motivated by the privacy applications discussed in the introduction, we introduce a parametrized class of decision problems in which the agent aims to identify *suitable* candidates—be it medical patients, investment opportunities, and so on—under privacy constraints. The agent is provided with a connection graph, such as in Figure 3 (left), which captures some notion of relationship or similarity between the candidates (possibly based on social ties, physical characteristics, geography, etc.). Connected subsets of suitable candidates, called *subgroups*, are present in the graph unknown to the agent. An action in the corresponding decision problem consists of selecting one of the candidates, and the goal is to maximize the number of distinct candidates selected from the suitable subgroups. If a selected candidate is a member of a subgroup, then the agent observes this fact henceforth; otherwise, the agent forgets having chosen this candidate at all (since it is sensitive info to know the unsuitable ones). Figure 3 (right) depicts an example; it resembles the prominent “Battleship” game, except that here the agent is absentminded about cells selected in the past that did not hit a ship. The parameters in this benchmark problem class control (a) the number of rounds

the agent can select candidates for; (b) an underlying graph structure, such as a 2D grid, or Erdős-Rényi $G(n, p)$ and $G(n, m)$ graph models; (c) the subgroups in the graph in terms of quantity, sizes, shapes (lines, cycles, cliques, stars) as well as the procedure for distributing them in the graph; and (d) the immediate payoffs for selecting suitable candidates of different subgroups.

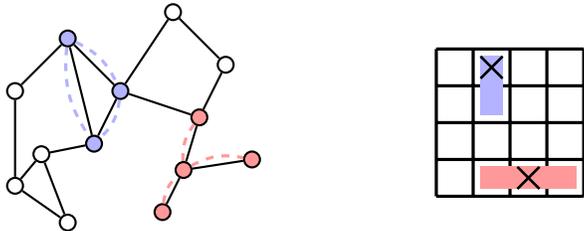


Figure 3: Subgroup detection under privacy constraints. On the left, we see an arbitrary graph with two subgroups (a 3-clique, and a star of degree 3). The goal is to find as many of the subgroups’ nodes as possible. On the right, we see another such decision problem on a 2D grid, which we visualized as an instance of the Absentminded Battleship game. The agent has already succeeded in hitting one node of each ship, which indicates that there must be more subgroup nodes nearby. The agent does not remember whether it has selected any cell other than these two before.

4.3 RANDOM DECISION PROBLEMS

Finally, we introduce a highly parametrized class of randomly generated decision problems, following an active line of work on random games (McKelvey & McLennan, 1996; Nudelman et al., 2004; Arieli & Babichenko, 2016; Amiet et al., 2021; Flesch et al., 2023; Heinrich et al., 2023). Part of their appeal is that they serve as a sanity check and help counterbalance cherry picking of benchmark problems. The parameters dictate (a) the probability with which a node will be terminal (dependent on its depth), (b) the probabilities with which a nonterminal node has k available actions, as well as with which it will be a chance node, (c) the (approximate) number of nodes we want to cover with each infoset, and (d) the probability distribution over payoffs at terminal nodes. The payoffs at the leaf nodes are drawn uniformly at random between 0 and 1. In the experiments in Section 5, each tree has varying depth in an interval $[d, d']$ where $4 \leq d \leq d' \leq 15$, the nonterminal nodes have 3 to 5 available actions and a 20% probability to be a chance node, and infosets are of a size roughly proportional to $n^{2/3}$, where n denotes the total number of decision nodes in the tree.

5 EXPERIMENTAL EVALUATION

Having introduced our benchmarks, we now use them to evaluate the performance of the algorithms described earlier in Section 3. Abbreviations “Sim,” “Det,” and “Rand” stand for simulation problems (Section 4.1), subgroup detection problems (Section 4.2), and random problems (Section 4.3) respectively. The suffixes indicate the number of nodes in the decision tree (with “k” and “m” abbreviating thousands and millions). Our algorithms run until any of three termination conditions is met: achieving a KKT gap of at most 10^{-6} , reaching the time limit of 4 hours,² or reaching the iteration limit of 6000. We run the first-order methods for 12 times with randomly initialized strategies and report the median. For GD and OGD, we run the algorithm with different learning rates, namely $\eta \in \{1, 10^{-1}, 10^{-2}, 10^{-3}\}$, and report only the one that minimizes the KKT gap the fastest at time of termination. We operate analogously for AMS, except that we instead test the parameter settings $(\eta, \beta_1, \beta_2) \in \{10^{-1}, 10^{-2}\} \times \{0.9, 0.99\} \times \{0.99, 0.999\}$ (this includes Reddi et al.’s suggested β -values). A subset of our results are gathered in Table 1 and aggregate statistics in Table 2. We also plot the KKT gap and value versus iteration in Figure 4 to gain insight into the process of convergence.³ The shaded regions represent the best and worst run for the respective iteration count and across the 12 initializations. Further experimental details and results can be found in the appendix.

²With the only exceptions of Det- $\{9m, 10m, 18m\}$ problems, which we run for 12 hours since the standard time limit poses a significant bottleneck for those instances.

³Our regret matching implementations complete more iterations per time than our gradient descent implementations, so the fact that we plot against iterations rather than time favors the gradient descent algorithms.

Problem	Gurobi			GD			OGD			AMS			RM			RM ⁺			PRM ⁺		
	value	time	gap	value	time	gap	value	time	gap	value	time	gap	value	time	gap	value	time	gap	value	time	gap
Det-1k	13.00	1m 24s	—	13.00	0.13s	—	13.00	0.07s	—	13.00	1.04s	—	13.00	0.32s	—	13.00	0.36s	—	13.00	0.41s	—
Det-1.8k	22.00	2m 40s	—	22.00	0.06s	—	22.00	0.07s	—	22.00	0.71s	—	22.00	0.03s	—	22.00	0.03s	—	22.00	0.03s	—
Det-2.0k	17.50	1m 42s	—	17.50	0.03s	—	17.50	0.05s	—	17.50	0.20s	—	17.50	0.03s	—	17.50	0.03s	—	17.50	0.03s	—
Det-2.1m	—	—	26.00	—	1e-05	25.96	—	0.02	26.15	—	0.002	26.15	—	0.003	26.15	3h 25m	—	26.15	—	0.005	
Det-2.2m	—	—	16.20	—	0.002	15.93	—	0.02	16.36	—	0.0002	16.36	2h 22m	—	16.36	3h 13m	—	16.36	—	5e-06	
Det-3.8m	—	—	15.66	—	0.003	15.14	—	0.03	15.78	—	0.0002	15.80	—	2e-06	15.80	—	5e-05	15.80	—	0.0003	
Det-9m	—	—	23.16	—	0.004	22.71	—	0.02	23.45	—	0.004	23.45	—	0.0001	23.45	—	0.0001	23.45	—	0.0004	
Det-10m	—	—	24.64	—	0.002	24.61	—	0.003	24.76	—	0.009	24.76	—	0.002	24.76	—	0.0004	24.76	—	0.0008	
Det-18m	—	—	26.38	—	0.006	25.81	—	0.05	26.71	—	0.004	26.71	—	0.004	26.71	—	0.001	26.71	—	0.04	
Rand-24k	0.72	—	0.66	7m 0s	—	0.66	7m 46s	—	0.66	4m 4s	—	0.66	26.55s	—	0.66	1m 3s	—	0.66	5m 5s	—	
Rand-35k	1.00	—	0.95	3.85s	—	0.95	3.76s	—	0.95	3.90s	—	0.92	0.99s	—	0.92	1.18s	—	0.94	1.68s	—	
Rand-42k	0.69	—	0.55	—	0.01	0.55	—	0.01	0.64	—	0.0006	0.65	—	2e-06	0.65	5m 56s	—	0.65	3m 19s	—	
Rand-13m	—	—	0.59	—	0.003	0.58	—	0.003	0.65	1h 40m	—	0.63	19m 11s	—	0.64	17m 31s	—	0.65	36m 42s	—	
Rand-18m	—	—	0.97	2h 33m	—	0.97	3h 0m	—	0.99	1h 29m	—	0.95	29m 45s	—	0.97	24m 0s	—	0.97	14m 31s	—	
Rand-23m	—	—	0.94	3h 37m	—	0.93	—	0.0007	0.98	3h 20m	—	0.98	23m 10s	—	0.96	23m 5s	—	0.95	18m 2s	—	
Sim-3k	6.25	1m 1s	—	6.25	0.32s	—	6.25	1.03s	—	6.25	5.54s	—	6.25	0.26s	—	6.25	0.28s	—	6.25	0.48s	—
Sim-7k	8.58	1m 36s	—	8.58	0.05s	—	8.58	0.05s	—	8.58	0.12s	—	8.58	0.05s	—	8.58	0.05s	—	8.58	0.05s	—
Sim-13k	10.38	4m 21s	—	10.38	0.69s	—	10.38	8.54s	—	10.38	14.37s	—	10.38	1.03s	—	10.38	1.01s	—	10.38	3.97s	—
Sim-540k	6.41	—	8.54	47.54s	—	8.54	2m 37s	—	8.54	15m 31s	—	8.54	19.39s	—	8.54	19.44s	—	8.54	3m 3s	—	
Sim-1m	4.14	—	4.77	5m 33s	—	4.77	7m 2s	—	4.77	29m 22s	—	4.77	2m 14s	—	4.77	2m 34s	—	4.77	4m 20s	—	
Sim-1.9m	—	—	13.45	18.31s	—	13.45	17.96s	—	13.45	1m 2s	—	13.45	12.36s	—	13.45	12.19s	—	13.45	12.47s	—	
Sim-2.3m	—	—	11.09	22.01s	—	11.09	21.88s	—	11.09	1m 10s	—	11.09	14.97s	—	11.09	15.00s	—	11.09	15.13s	—	
Sim-4m	—	—	14.01	45m 5s	—	14.01	41m 0s	—	13.98	—	0.02	14.01	11m 36s	—	14.01	7m 3s	—	14.01	21m 17s	—	

Table 1: The performance of various algorithms in a subset of our benchmarks. Value and convergence winners per game highlighted in bold. For Gurobi, time is only reported if it terminated.

Category	Metric	Gurobi	GD	OGD	AMS	RM	RM ⁺	PRM	PRM ⁺
Utility value	% of best (↑)	36.1	42.6	42.6	78.7	65.6	70.5	72.1	72.1
	avg. rank (↓)	6.26	5.12	5.34	3.60	4.04	3.83	3.89	3.91
Convergence	% reached (↑)	27.9	72.1	72.1	77.0	83.6	90.2	78.7	80.3
	% of best (↑)	0.0	11.5	6.6	0.0	41.0	39.3	23.0	21.3
	avg. rank (↓)	8.00	4.80	5.52	5.74	2.34	2.17	3.84	3.59

Table 2: Aggregate statistics of the full experiments in the appendix. We report how often an algorithm reached convergence across all 61 benchmark instances, and relative to the other algorithms, how well it ranks on average (lower is better), and how often it achieved the best value / convergence.

The main takeaways are the following:

- Gurobi fails to converge beyond small instances ($\leq 100k$ nodes for simulation, and $\leq 20k$ otherwise). Moreover, when it reaches termination, the time required is multiple orders of magnitude more than that of the first-order optimizers. This is despite the fact that Gurobi is based on an optimized C++ implementation whereas our first-order optimizers are implemented in Python.
- Interestingly, in all such cases in Table 1, where we know the optimal value, the first-order optimizers converge to an optimal strategy.

As expected, we can also find some experiments where the previous point will not hold (e.g. Rand-42k once Gurobi would eventually terminate). Indeed, we design an extreme example in the appendix where our gradient descent and regret matching methods all converge to a KKT point that is arbitrarily bad in value relative to the global optimum. We also illustrate on two examples that one of the two algorithm types can perform arbitrarily bad while the other reaches global optimum.

- AMS and RM⁺ oftentimes attain higher values than GD and OGD, and almost never less.
- The RM family of algorithms, and RM⁺ in particular, consistently outperform GD, OGD, and AMS in runtime. The difference is often many orders of magnitude, especially in the larger instances.

The literature on two-player zero-sum *perfect-recall* games may offer some conceptual explanation here. Namely, it is believed that the following property of RM methods—which continues to be satisfied in our setting with imperfect recall—is important to practical performance: at the infostep level, RM methods are invariant to step sizes (Chakrabarti et al., 2024). Instead, they regulate their step sizes based on the current and past gradients alone, requiring no careful step size tuning in the same way that gradient descent does.

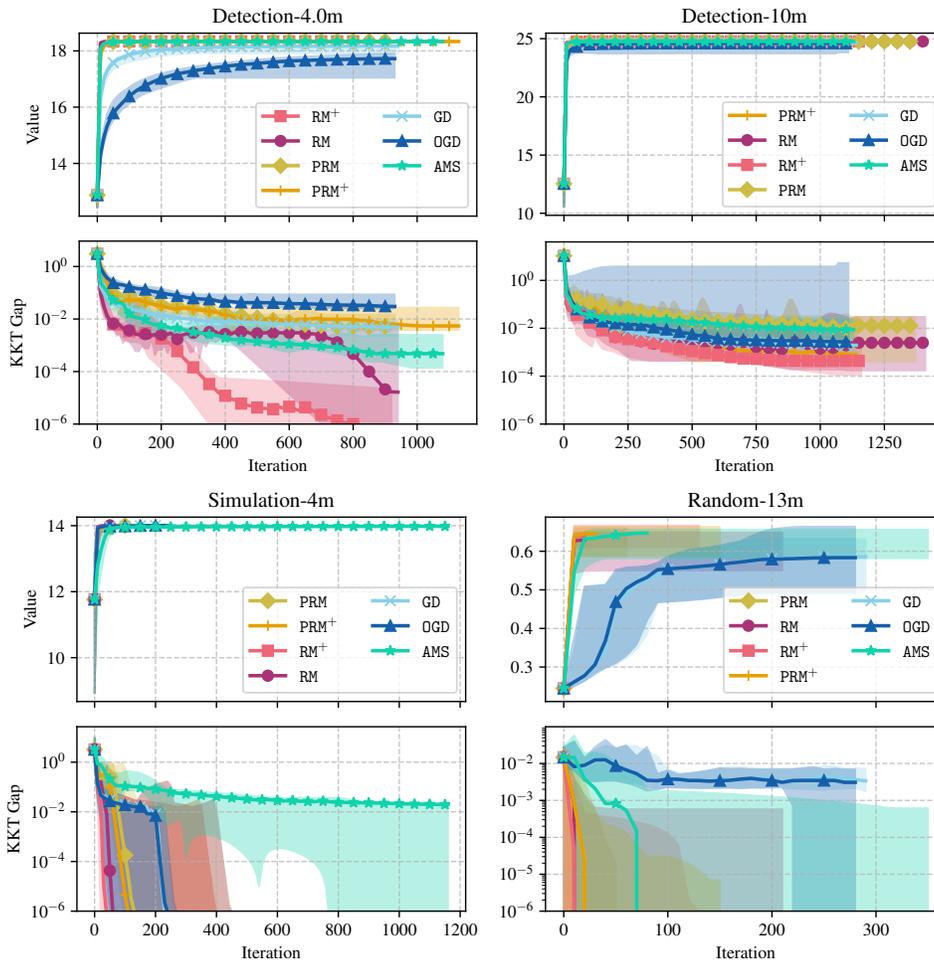


Figure 4: Two detection instances, a simulation and a random instance. They have ~ 1200 infosets, ~ 300 infosets, 3 infosets, and ~ 100 infosets respectively.

5. RM^+ performs best among the RM family. Surprisingly, it typically outperforms PRM^+ , which stands in stark contrast to what has been observed in zero-sum games (Farina et al., 2021).

Our intuition for that is as follows. Predictiveness in RM (= Optimism in GD) roughly corresponds to having negative momentum, which is beneficial in zero-sum games and minimax optimization because it helps minimize regret faster. But in our setting of nonlinear (single-player) optimization, it is not known whether predictiveness helps anymore, since the task is not to minimize regret but to search for a first-order optimal point. Indeed, our experiments seem to suggest otherwise.

6 FUTURE RESEARCH

Our paper opens many interesting avenues for future work. First, we have focused exclusively on solving tabular imperfect-recall decision problems. A promising direction is to use modern RL techniques to expand the scope to even larger problems that cannot be represented in tabular form. Considering other formulations beyond tree-form decision problems, such as (PO)MDPs, is another natural direction that was beyond our scope. Finally, our experiments revealed that the regret matching family of algorithms is a formidable first-order optimizer; elucidating their theoretical properties is another important open question.

540 REPRODUCIBILITY STATEMENT

541
542 Comprehensive details about our experimental methodology can be found in Sections 4 and 5 as
543 well as in the last section of the appendix, including generation procedures, hyperparameter ranges,
544 termination conditions, hyperparameter grids, and hardware specifications. The supplementary ma-
545 terial contains our code base for generating instances from our benchmark suite, running the dis-
546 cussed algorithms on problem instances, and evaluating and visualizing the results. The specific 61
547 benchmark instances used in our evaluation, along with their corresponding experimental results,
548 are available through the provided link at the end of the appendix.

549
550 REFERENCES

551
552 Ben Amiet, Andrea Collevocchio, Marco Scarsini, and Ziwen Zhong. Pure nash equilibria and best-
553 response dynamics in random games. *Mathematics of Operations Research*, 46(4):1552–1572,
554 2021.

555 Itai Arieli and Yakov Babichenko. Random extensive form games. *Journal of Economic Theory*,
556 166:517–535, 2016.

557
558 Kenneth Arrow. Economic welfare and the allocation of resources for invention. In *The Rate and*
559 *Direction of Inventive Activity: Economic and Social Factors*, pp. 609–626. National Bureau of
560 Economic Research, Inc, 1962.

561 Nicola Basilico, Andrea Celli, Giuseppe De Nittis, and Nicola Gatti. Team-maxmin equilibrium:
562 Efficiency bounds and algorithms. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2017.

563
564 Heymann Benjamin and Marc Lanctot. Learning in games with progressive hiding.
565 *arXiv:2409.03875*, 2024.

566
567 Ratip Emin Berker, Emanuel Tewelde, Ioannis Anagnostides, Tuomas Sandholm, and Vincent
568 Conitzer. The value of recall in extensive-form games. In *AAAI Conference on Artificial In-*
569 *telligence (AAAI)*, 2025.

570 David Blackwell. An analog of the minmax theorem for vector payoffs. *Pacific Journal of Mathe-*
571 *matics*, 6:1–8, 1956.

572
573 Michael Bowling, Neil Burch, Michael Johanson, and Oskari Tammelin. Heads-up limit hold'em
574 poker is solved. *Science*, 347(6218):145–149, 2015.

575 Rachael Briggs. Putting a value on beauty. In Tamar Szabó Gendler and John Hawthorne (eds.),
576 *Oxford Studies in Epistemology: Volume 3*, pp. 3–34. Oxford University Press, 2010.

577
578 Noam Brown and Tuomas Sandholm. Superhuman AI for heads-up no-limit poker: Libratus beats
579 top professionals. *Science*, 359(6374):418–424, 2018.

580
581 Noam Brown and Tuomas Sandholm. Superhuman AI for multiplayer poker. *Science*, 365(6456):
582 885–890, 2019.

583 Jiří Čermák, Branislav Bosanský, and Viliam Lisý. An algorithm for constructing and solving
584 imperfect recall abstractions of large extensive-form games. In *International Joint Conference on*
585 *Artificial Intelligence (IJCAI)*, 2017a.

586
587 Jiří Čermák, Branislav Bošanský, and Michal Pěchouček. Combining incremental strategy gener-
588 ation and branch and bound search for computing maxmin strategies in imperfect recall games.
589 In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems, AAMAS*
590 *'17*, pp. 902–910. International Foundation for Autonomous Agents and Multiagent Systems,
591 2017b.

592 Jiří Čermák, Branislav Bošanský, Karel Horák, Viliam Lisý, and Michal Pěchouček. Approximating
593 maxmin strategies in imperfect recall games using a-loss recall property. *International Journal of*
Approximate Reasoning, 93:290–326, 2018.

- 594 Jiří Čermák, Viliam Lisý, and Branislav Bošanský. Automated construction of bounded-loss
595 imperfect-recall abstractions in extensive-form games. *Artificial Intelligence*, 282, 2020.
596
- 597 Darshan Chakrabarti, Julien Grand-Clément, and Christian Kroer. Extensive-form game solving via
598 blackwell approachability on treeplexes. In *Advances in Neural Information Processing Systems*
599 38, 2024.
- 600 Eric O. Chen, Alexis Ghersengorin, and Sami Petersen. Imperfect recall and AI delegation, 2024.
601
- 602 Emery Cooper, Caspar Oesterheld, and Vincent Conitzer. Characterising simulation-based program
603 equilibria. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2025.
- 604 Constantinos Daskalakis and Ioannis Panageas. The limit points of (optimistic) gradient descent in
605 min-max optimization. In *NeurIPS 2018*, pp. 9256–9266, 2018.
606
- 607 Constantinos Daskalakis and Christos Papadimitriou. Continuous local search. In *ACM-SIAM Sym-*
608 *posium on Discrete Algorithms (SODA)*, 2011.
- 609 Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training gans with
610 optimism. In *International Conference on Learning Representations (ICLR)*, 2018.
611
- 612 Adam Elga. Self-locating belief and the Sleeping Beauty problem. *Analysis*, 60(2):143–147, 2000.
613
- 614 Scott Emmons, Caspar Oesterheld, Andrew Critch, Vincent Conitzer, and Stuart Russell. For learn-
615 ing in symmetric teams, local optima are global nash equilibria. In *International Conference on*
616 *Machine Learning (ICML)*, 2022.
- 617 European Parliament and Council of the EU. Regulation (EU) 2016/679 (General Data Protection
618 Regulation), 2016. URL <https://eur-lex.europa.eu/eli/reg/2016/679/oj>. Of-
619 ficial Journal of the European Union.
- 620 Gabriele Farina, Andrea Celli, Nicola Gatti, and Tuomas Sandholm. Ex ante coordination and collu-
621 sion in zero-sum multi-player extensive-form games. In *Neural Information Processing Systems*
622 *(NeurIPS)*, 2018.
623
- 624 Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Faster game solving via predictive Black-
625 well approachability: Connecting regret matching and mirror descent. In *AAAI Conference on*
626 *Artificial Intelligence (AAAI)*, 2021.
- 627 Gabriele Farina, Julien Grand-Clément, Christian Kroer, Chung-Wei Lee, and Haipeng Luo. Re-
628 regret matching+: (in)stability and fast convergence in games. In *Neural Information Processing*
629 *Systems (NeurIPS)*, 2023.
630
- 631 John Fearnley, Paul Goldberg, Alexandros Hollender, and Rahul Savani. The complexity of gradient
632 descent: $\text{CLS} = \text{PPAD} \cap \text{PLS}$. *Journal of the ACM*, 70(1):7:1–7:74, 2023.
- 633 János Flesch, Arkadi Predtetchinski, and Ville Suomala. Random perfect information games. *Math-*
634 *ematics of Operations Research*, 48(2):708–727, 2023.
635
- 636 Dean P. Foster and Sergiu Hart. Smooth calibration, leaky forecasts, finite recall, and nash dynamics.
637 *Games and Economic Behavior*, 109:271–293, 2018.
- 638 Drew Fudenberg and Jean Tirole. *Game Theory*. MIT Press, 1991.
639
- 640 Hugo Gimbert, Soumyajit Paul, and B. Srivathsan. A bridge between polynomial optimization and
641 games with imperfect recall. In *Autonomous Agents and Multi-Agent Systems*, 2020.
- 642 Hugo Gimbert, Soumyajit Paul, and B. Srivathsan. Simplifying imperfect recall games. In *Pro-*
643 *ceedings of the 24th International Conference on Autonomous Agents and Multiagent Systems*,
644 AAMAS ’25, pp. 895–903. International Foundation for Autonomous Agents and Multiagent
645 Systems, 2025.
646
- 647 Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2025. URL <https://www.gurobi.com>.

- 648 Sergiu Hart and Andreu Mas-Colell. A simple adaptive procedure leading to correlated equilibrium.
649 *Econometrica*, 68:1127–1150, 2000.
650
- 651 Sergiu Hart and Andreu Mas-Colell. Regret-based continuous-time dynamics. *Games and Economic*
652 *Behavior*, 45(2):375–394, 2003.
653
- 654 Torsten Heinrich, Yoojin Jang, Luca Mungo, Marco Pangallo, Alex Scott, Bassel Tarbush, and
655 Samuel Wiese. Best-response dynamics, playing sequences, and convergence to equilibrium in
656 random games. *International Journal of Game Theory*, 52:703–735, 2023.
- 657 Richard Helgason, Jeffery Kennington, and H Lall. A polynomially bounded algorithm for a singly
658 constrained quadratic program. *Mathematical Programming*, 18:338–343, 1980.
659
- 660 J. R. Isbell. *Finitary Games*, pp. 79–96. Princeton University Press, 1957.
661
- 662 Mamoru Kaneko and J Jude Kline. Behavior strategies, mixed strategies and perfect recall. *Inter-*
663 *national Journal of Game Theory*, 24:127–145, 1995.
- 664 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd Interna-*
665 *tional Conference on Learning Representations, ICLR 2015*, 2015.
666
- 667 Megan Kinniment, Lucas Jun Koba Sato, Haoxing Du, Brian Goodrich, Max Hasin, Lawrence Chan,
668 Luke Harold Miles, Tao R. Lin, Hjalmar Wijk, Joel Burget, Aaron Ho, Elizabeth Barnes, and Paul
669 Christiano. Evaluating language-model agents on realistic autonomous tasks. *arXiv:2312.11671*,
670 2024.
- 671 J Jude Kline. Minimum memory for equivalence between ex ante optimality and time-consistency.
672 *Games and Economic Behavior*, 38(2):278–305, 2002.
673
- 674 Daphne Koller and Nimrod Megiddo. The complexity of two-person zero-sum games in extensive
675 form. *Games and Economic Behavior*, 4(4):528–552, 1992.
- 676 Vojtěch Kovařík, Caspar Oesterheld, and Vincent Conitzer. Game theory with simulation of other
677 players. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2023.
678
- 679 Vojtěch Kovařík, Caspar Oesterheld, and Vincent Conitzer. Recursive joint simulation in games.
680 *arXiv:2402.08128*, 2024.
- 681 Vojtěch Kovařík, Eric Olav Chen, Sami Petersen, Alexis Ghersengorin, and Vincent Conitzer. Ai
682 testing should account for sophisticated strategic behaviour. In *Neural Information Processing*
683 *Systems (NeurIPS)*, 2025a.
684
- 685 Vojtěch Kovařík, Nathaniel Sauerberg, Lewis Hammond, and Vincent Conitzer. Game theory with
686 simulation in the presence of unpredictable randomisation. In *Autonomous Agents and Multi-*
687 *Agent Systems*, 2025b.
- 688
- 689 Christian Kroer and Tuomas Sandholm. Extensive-form game abstraction with bounds. In *ACM*
690 *Conference on Economics and Computation (EC)*, 2014.
- 691
- 692 Christian Kroer and Tuomas Sandholm. Imperfect-recall abstractions with bounds in games. In
693 *ACM Conference on Economics and Computation (EC)*, 2016.
- 694
- 695 H. W. Kuhn. Extensive games and the problem of information. In *Contributions to the Theory of*
696 *Games*, volume 2 of *Annals of Mathematics Studies*, 28, pp. 193–216. Princeton University Press,
1953.
- 697
- 698 Nicolas S Lambert, Adrian Marple, and Yoav Shoham. On equilibria in games with imperfect recall.
699 *Games and Economic Behavior*, 113:164–185, 2019.
- 700
- 701 Marc Lanctot, Richard Gibson, Neil Burch, Martin Zinkevich, and Michael Bowling. No-regret
learning in extensive-form games with imperfect recall. In *International Conference on Machine*
Learning (ICML), 2012.

- 702 Boning Li and Longbo Huang. Efficient online pruning and abstraction for imperfect information
703 extensive-form games. In *The Thirteenth International Conference on Learning Representations,*
704 *ICLR 2025*. OpenReview.net, 2025.
- 705
- 706 Mingyang Liu, Gabriele Farina, and Asuman Ozdaglar. LiteEFG: An efficient python library for
707 solving extensive-form games. *arXiv:2407.20351*, 2024.
- 708
- 709 Tai-Yu Ma and Philippe Gerber. Distributed regret matching algorithm for dynamic congestion
710 games with information provision. *Transportation Research Procedia*, 3:3–12, 2014.
- 711
- 712 Richard D. McKelvey and Andrew McLennan. Computation of equilibria in finite games. In
713 H. Amann, D. Kendrick, and J. Rust (eds.), *Handbook of Computational Economics*, volume 1,
714 pp. 87–142. Elsevier, 1996.
- 715
- 716 Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor
717 Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial
718 intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017.
- 719
- 720 Richard R Nelson. The simple economics of basic scientific research. *Journal of political economy*,
721 67(3):297–306, 1959.
- 722
- 723 Eugene Nudelman, Jennifer Wortman, Kevin Leyton-Brown, and Yoav Shoham. Run the GAMUT:
724 A comprehensive approach to evaluating game-theoretic algorithms. In *International Conference*
725 *on Autonomous Agents and Multi-Agent Systems (AAMAS)*, New York, NY, 2004.
- 726
- 727 Caspar Oesterheld. Robust program equilibrium. *Theory and Decision*, 86(1):143–159, 2 2019.
- 728
- 729 Caspar Oesterheld and Vincent Conitzer. Can *de se* choice be *ex ante* reasonable in games of imper-
730 fect recall? a complete analysis. [https://www.andrew.cmu.edu/user/coesterh/](https://www.andrew.cmu.edu/user/coesterh/DeSeVsExAnte.pdf)
731 [DeSeVsExAnte.pdf](https://www.andrew.cmu.edu/user/coesterh/DeSeVsExAnte.pdf), 2024. Working paper. Accessed: 2024-07-13.
- 732
- 733 Alexander Pan, Jun Shern Chan, Andy Zou, Nathaniel Li, Steven Basart, Thomas Woodside, Han-
734 lin Zhang, Scott Emmons, and Dan Hendrycks. Do the rewards justify the means? measuring
735 trade-offs between rewards and ethical behavior in the machiavelli benchmark. In *International*
736 *Conference on Machine Learning (ICML)*, 2023.
- 737
- 738 Christos H. Papadimitriou and Mihalis Yannakakis. On complexity as bounded rationality (extended
739 abstract). In *Symposium on Theory of Computing (STOC)*, 1994.
- 740
- 741 Michele Piccione and Ariel Rubinstein. On the interpretation of decision problems with imperfect
742 recall. *Games and Economic Behavior*, 20:3–24, 1997.
- 743
- 744 L.D. Popov. A modification to the Arrow-Hurwicz method for search of saddle-points. *Mathemati-*
745 *cal Notes of the Academy of Sciences of the USSR*, 28(5):845–848, 1980.
- 746
- 747 Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *6th*
748 *International Conference on Learning Representations, ICLR*, 2018.
- 749
- 750 Ariel Rubinstein. *Modeling Bounded Rationality*. MIT Press, 1998.
- 751
- 752 Iosif Sakos, Stefanos Leonardos, Stelios Andrew Stavroulakis, Will Overman, Ioannis Panageas,
753 and Georgios Piliouras. Beating price of anarchy and gradient descent without regret in potential
754 games. In *International Conference on Learning Representations (ICLR)*, 2024.
- 755
- 756 Matthew Stephenson, Andrew Miller, Xyn Sun, Bhargav Annem, and Rohan Parikh. NDAI agree-
757 ments. *arXiv:2502.07924*, 2025.
- 758
- 759 Moshe Tennenholtz. Program equilibrium. *Games and Economic Behavior*, 49(2):363–373, 2004.
- 760
- 761 Emanuel Tewolde, Caspar Oesterheld, Vincent Conitzer, and Paul W. Goldberg. The computational
762 complexity of single-player imperfect-recall games. In *International Joint Conference on Artificial*
763 *Intelligence (IJCAI)*, 2023.

- 756 Emanuel Tewolde, Brian Hu Zhang, Caspar Oesterheld, Manolis Zampetakis, Tuomas Sandholm,
757 Paul W. Goldberg, and Vincent Conitzer. Imperfect-recall games: Equilibrium concepts and their
758 complexity. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2024.
759
- 760 Emanuel Tewolde, Brian Hu Zhang, Caspar Oesterheld, Tuomas Sandholm, and Vincent Conitzer.
761 Computing game symmetries and equilibria that respect them. In *AAAI Conference on Artificial
762 Intelligence (AAAI)*, 2025.
- 763 Bernhard Von Stengel and Daphne Koller. Team-maxmin equilibria. *Games and Economic Behav-*
764 *ior*, 21(1-2):309–321, 1997.
765
- 766 Kevin Waugh. Abstraction in large extensive games. Master’s thesis, University of Alberta, 2009.
- 767 Kevin Waugh, Martin Zinkevich, Michael Johanson, Morgan Kan, David Schnizlein, and Michael
768 Bowling. A practical use of imperfect recall. In *Symposium on Abstraction, Reformulation and
769 Approximation (SARA)*, 2009.
770
- 771 Chen-Yu Wei, Chung-Wei Lee, Mengxiao Zhang, and Haipeng Luo. Linear last-iterate convergence
772 in constrained saddle-point optimization. In *International Conference on Learning Representa-*
773 *tions (ICLR)*, 2021.
- 774 Brian Hu Zhang and Tuomas Sandholm. Polynomial-time optimal equilibria with a mediator in
775 extensive-form games. In *Neural Information Processing Systems (NeurIPS)*, 2022.
776
- 777 Brian Hu Zhang, Gabriele Farina, and Tuomas Sandholm. Team belief DAG: generalizing the se-
778 quence form to team games for fast computation of correlated team max-min equilibria via regret
779 minimization. In *International Conference on Machine Learning (ICML)*, 2023.
- 780 Youzhi Zhang, Bo An, and V. S. Subrahmanian. Correlation-based algorithm for team-maxmin
781 equilibrium in multiplayer extensive-form games. In *International Joint Conference on Artificial
782 Intelligence (IJCAI)*, 2022.
- 783 Martin Zinkevich, Michael Bowling, Michael Johanson, and Carmelo Piccione. Regret minimization
784 in games with incomplete information. In *Neural Information Processing Systems (NIPS)*, 2007.
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810 A APPENDIX

811 Here, we expand further on details we omitted in the main body.

814 A.1 FURTHER RELATED WORK

815 **Simulation games** Commencing from the paper of Kovařík et al. (2023), there has been significant
816 interest in situations where one player can simulate another player (Chen et al., 2024; Kovařík et al.,
817 2024; 2025b; Oesterheld, 2019; Cooper et al., 2025); this is precisely the type of problem captured
818 by one of our benchmarks. The premise of simulating the other player is strongly connected with
819 the notion of *program equilibrium* (Tennenholtz, 2004), where players are allowed to submit source
820 code. This turns out to unlock more cooperative outcomes by expanding the set of equilibria.
821

822 **MDPs and repeated games** Another notable motivation for examining imperfect-recall decision
823 problems lies in the fact that they can result in simpler and more interpretable strategies. This
824 point can be illustrated well in the context of Markov decision problems (MDPs), where insisting
825 on *Markovian* policies—which depend solely on the state and not the entire history—is particu-
826 larly common; this can be viewed as an extreme form of imperfect recall. Relatedly, restricting
827 the memory and description complexity of a policy has received a lot of attention in the context
828 of repeated games (e.g., Foster & Hart, 2018; Papadimitriou & Yannakakis, 1994). In certain set-
829 tings, near-optimal policies are possible even under imperfect recall. More broadly, the question of
830 characterizing the value of recall was recently addressed by Berker et al. (2025).
831

832 **CDT equilibria** The CDT equilibrium falls under the family based on the *multi-selves* ap-
833 proach (Kuhn, 1953). At a high level, whenever the player has imperfect information on the decision
834 node of an infoset it is currently in, the player will weight each possibility with the probability of
835 reaching the decision node in question under strategy x . The name is derived from the intuition
836 that the player’s choice to deviate from x at the current node does not *cause* any change in its
837 behavior at any other node, even if they are of the same infoset. Another prominent member is
838 the *EDT equilibrium*, which results from marrying evidential decision theory with the multi-selves
839 approach (Oesterheld & Conitzer, 2024). For further background on the ongoing debate around
840 decision theories and how they relate to belief formation (*cf.* the “sleeping beauty” problem (Elga,
841 2000)), we refer to (Piccione & Rubinstein, 1997; Briggs, 2010; Oesterheld & Conitzer, 2024). Fur-
842 ther, we refer to Tewolde et al. (2024) for a computational treatment of equilibria in multi-player
843 games with imperfect recall. With regard to the complexity of computing CDT equilibria, we saw
844 earlier that a $\text{poly}(1/\epsilon)$ time algorithm exists by running GD on a suitable optimization problem; in
845 the regime where ϵ is exponentially small, the complexity is characterized by the class CLS, and
846 is believed to be hard (Daskalakis & Papadimitriou, 2011; Fearnley et al., 2023; Tewolde et al.,
847 2023). Conceptually, and also computationally, CDT equilibria in decision problems with imperfect
848 recall have been also connected to Nash equilibria in team games that respect a given set of game
849 symmetries (Lambert et al., 2019; Emmons et al., 2022; Tewolde et al., 2025).

850 **Regret matching** Regret matching and its variants have received a lot of attention in (two-player)
851 zero-sum *perfect recall* extensive-form games. In particular, the *counterfactual regret minimization*
852 (*CFR*) algorithm, famously introduced by Zinkevich et al. (2007), employs a separate RM algorithm
853 for each information set. The CFR framework has spawned a flourishing, and still active, line of
854 work. Yet, much less is known beyond (two-player) zero-sum games. It has to be stressed again that
855 in zero-sum games, RM and its variants only have guarantees concerning the time average strategy. In
856 fact, the last iterate can fail to converge (Farina et al., 2023). Our experiments suggest a fundamental
857 difference in constrained optimization problems: all our results make use of the last iterate, which
858 not only converges, but does so remarkably fast. To our knowledge, there is currently no theory that
859 predicts that RM and its variants will converge. The continuous time of RM was analyzed by Hart
860 & Mas-Colell (2003), who also established asymptotic convergence in two-player potential games
861 for a certain—somewhat artificial—variant of RM in discrete time. Fast empirical convergence was
862 reported by Ma & Gerber (2014) in a certain class of congestion games.

863 An intriguing behavior we uncover in this paper is that the RM family of algorithms often outperforms
(O)GD in terms of the attained value, at least for the benchmark problems we consider. In the context
of multi-player potential games, which is closely related to imperfect-recall decision problems, the

864 problem of characterizing the performance of different algorithms is poorly understood. One notable
 865 contribution here is the recent paper of Sakos et al. (2024), but it only focused on 2×2 games.
 866 Providing a theoretical explanation that justifies the excellent performance of RM in terms of value
 867 is an interesting but challenging direction for the future.

868
 869 **Game abstractions and related work** As stated in the introduction, games with imperfect recall
 870 have found great success in the state-of-the-art algorithms for solving real-world games that are too
 871 large to handle, and therefore need to be compressed in a game abstraction (Waugh et al., 2009;
 872 Čermák et al., 2017a; Brown & Sandholm, 2018; 2019; Benjamin & Lanctot, 2024; Li & Huang,
 873 2025). The motivation and techniques from that line of literature are complementary to ours because

- 874 1. they can flexibly change the forms and structures of imperfect recall since there is an
 875 underlying *perfect-recall* game that the user cares about and an imperfect-recall abstraction
 876 is evaluated by how well its solutions can be lifted to good strategies in the underlying
 877 perfect-recall game, and because
- 878 2. their forms of imperfect recall are *designed* to be relatively benign so that the solutions in
 879 the imperfect-recall abstraction can be computed efficiently (recall the hardness of Propo-
 880 sition 2).

882 Our work, on the other hand, is interested in decision problems that inherently exhibit imperfect
 883 recall, and solving them for their own sake (there is no underlying perfect-recall problem), even if
 884 they are computationally hard (such as when absentminded infosets are present).

885 Various narrower forms of imperfect recall have received wide attention. Decision problems without
 886 absentmindedness always admit optimal solutions in pure strategies, which enables methods based
 887 on mixed-integer programming and double oracle-style incremental strategy generation (Čermák
 888 et al., 2017a;b; 2020). Prominent forms of imperfect recall that are computationally benign include
 889 skew well-formed games, for which CFR is still guaranteed to work (Lanctot et al., 2012; Kroer &
 890 Sandholm, 2016), and decision problems with A-loss recall (Kaneko & Kline, 1995; Kline, 2002),
 891 for which optimal strategies can be computed in polynomial time (Čermák et al., 2018; Gimbert
 892 et al., 2025). CFR methods have also shown experimental success in slightly more general settings
 893 of imperfect-recall abstractions (Kroer & Sandholm, 2016; Čermák et al., 2020; Li & Huang, 2025),
 894 though these settings continue to be narrow subclasses of decision problems *without absentminded-*
 895 *ness*. Indeed, CFR is a framework designed for the perfect recall setting, and past work—such as
 896 Waugh et al. (2009, Section “Challenges of Imperfect Recall”)—have discussed why CFR concep-
 897 tually cannot be extended to imperfect-recall settings beyond narrow subclasses such as (variants
 898 of) skew well-formed games. CFR updates the action probabilities at each infoset based on a notion
 899 of expected utilities that counterfactually assumes that the player played in the past as if it only
 900 wanted to reach the infoset in question. Under imperfect recall, and especially absentmindedness,
 901 notions such as “in the past” and “playing actions in order to reach an infoset” (once, or multiple
 902 times?) become highly dubious. In order to apply RM methods on decision problems with arbitrary
 903 forms of imperfect recall, we step away from CFR in this paper, and instead work in the *agent-form*
 904 of the decision problem Kuhn (1953). The agent-form imagines each infoset as being played by a
 905 separate player, therefore introducing multi-agent equilibrium concepts even in single-agent deci-
 906 sion problems. Moreover, we tackle infosets with absentmindedness through causal decision theory
 907 (and a belief formation system that is compatible with it; *cf.* Tewolde et al., 2024). This equates
 908 to regret matching methods being applied on gradients from the polynomial optimization problem
 909 Equation (1). Last but not least, we are not studying RM methods as regret minimizers—as it is
 910 the case in the literature on imperfect-recall abstractions, using the CFR framework—but instead
 911 evaluate RM methods in their performance as a first-order optimizers. These two are incomparable.

912 **Mixed strategies and team games** Much of the prior work in extensive-form games has focused
 913 on *mixed* strategies—probability distributions over pure strategies. Unlike behavioral strategies,
 914 mixed strategies allow the player to correlate its actions across infosets; one such example is *ex-*
 915 *ante* team coordination (Farina et al., 2018) in the context of team games. As we explained in our
 916 introduction, a team game can be phrased as an imperfect-recall decision problem; in fact, one with-
 917 out absentmindedness. Without absentmindedness, it follows that there exists an optimal strategy
 that is pure; in contrast, the presence of absentmindedness—which is primary focus on this paper—
 requires randomization (Isbell, 1957). In the presence of imperfect recall, mixed strategies are not

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

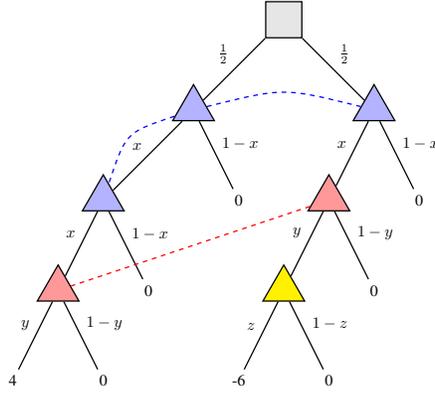


Figure 5: This is the decision making instance one would obtain from applying the construction of the proof of Theorem 1 to the polynomial maximization $\max 2x^2y - 3xyz$ s.t. $0 \leq x, y, z \leq 1$.

realization-equivalent to behavioral strategies (Kuhn, 1953), and they do not fit our motivation since they imply a form of memory mechanism. Related to *ex-ante* team coordination, classical equilibrium concepts in extensive-form games involving *correlation* can be modeled via a mediator—a trusted third party—with imperfect recall (Zhang & Sandholm, 2022); that the mediator has imperfect recall can serve to safeguard the players’ sensitive information, which is tied to one of the key motivations of this paper.

A.2 ON THEOREM 1: AN EQUIVALENCE TO CONSTRAINED POLYNOMIAL OPTIMIZATION

Starting with part 1 of Theorem 1, the concrete polynomial optimization problem for maximizing utility in a decision making instance takes the following form:

$$\begin{aligned} \max_{\mathbf{x} \in \times_{I \in \mathcal{I}} \mathbb{R}^{A_I}} \quad & U(\mathbf{x}) = \sum_{z \in \mathcal{Z}} u(z) \cdot \mathbb{P}(z \mid \mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x}(a \mid I) \geq 0 \quad \forall I \in \mathcal{I}, \forall a \in A_I \\ & \sum_{a \in A_I} \mathbf{x}(a \mid I) = 1 \quad \forall I \in \mathcal{I} \end{aligned} \quad (1)$$

Recall that $\mathbb{P}(z \mid \mathbf{x})$ is a product of action probabilities $\mathbf{x}(a \mid I)$ and chance probabilities of actions the agent and chance need to take in order to reach leaf node z . The former kind are the optimization variables in the above program, and the latter kind are fixed scalars.

Next, we move our attention to part 2 of Theorem 1. As an example, take the polynomial maximization instance

$$\max 2x^2y - 3xyz \quad \text{s.t.} \quad 0 \leq x, y, z \leq 1.$$

Figure 5 then depicts its corresponding decision making instance under imperfect recall that one obtains from the construction in the proof of Theorem 1. The more general construction idea is as follows. Variables correspond to infosets, and occurrences of a variable to a decision node. For a hypercube domain, each decision node will have two actions; for general products of simplices, the number of actions at a decision node associated to a variable x is precisely the number of vertices in the simplex that constrains x . A chance node at the root selects uniform randomly among the different monomials, and the utility payoffs are obtained from the coefficients of the monomials times the number of monomials present in the polynomial function. A nonzero utility payoff is only obtained if the player selected the left action for every variable occurrence in the monomial that was drawn by the chance node.

A.3 A NOTE ON OUR GUROBI IMPLEMENTATION

Our Gurobi method implements the constrained polynomial optimization problem (1). Gurobi supports nonlinear optimization if they can be implemented via quadratic constraints. Therefore, we

Algorithm 5: AMSGrad; AMS

```

1 Initialize learning rate  $\eta > 0$ ,  $\beta_1, \beta_2 \in [0, 1)$ ,  $\mathbf{x}^{(1)} \in \Delta(m)$ , and  $\mathbf{m}^{(0)}, \mathbf{v}^{(0)}, \hat{\mathbf{v}}^{(0)} = \mathbf{0}$ 
2 procedure GETX( $\tilde{\mathbf{u}}^{(t)}$ ) return  $\mathbf{x}^{(t)}$ 
3 procedure STEP( $\mathbf{u}^{(t)}$ )
4    $\mathbf{m}^{(t)} \leftarrow \beta_1 \mathbf{m}^{(t-1)} + (1 - \beta_1) \mathbf{u}^{(t)}$ 
5    $\mathbf{v}^{(t)} \leftarrow \beta_2 \mathbf{v}^{(t-1)} + (1 - \beta_2) (\mathbf{u}^{(t)})^2$ 
6    $\hat{\mathbf{v}}^{(t)} \leftarrow \max\{\hat{\mathbf{v}}^{(t-1)}, \mathbf{v}^{(t)}\}$ 
7    $\mathbf{x}^{(t+1)} \leftarrow \Pi_{\Delta(m), \sqrt{\hat{\mathbf{v}}^{(t)}}}(\mathbf{x}^{(t)} + \eta \mathbf{m}^{(t)} / \hat{\mathbf{v}}^{(t)})$ 

```

follow the common practice to implement a product of variables $\prod_{i=1}^k x_i$ in the objective function—which can arise from the terms $\mathbb{P}(z \mid \mathbf{x})$ —as follows. We introduce $k - 1$ auxiliary variables y_2, \dots, y_k , replace $\prod_{i=1}^k x_i$ in the objective with the single variable y_k , and add the following $k - 1$ quadratic equations as additional constraints to the program: $y_2 = x_2 \cdot x_1$, and for each $i = 3, \dots, k$: $y_i = x_i \cdot y_{i-1}$. If we work on a decision tree—as we do in this paper—we can reuse auxiliary variables efficiently in the following way. The same auxiliary variable y_i will appear for all paths in the tree from the root node h_0 to some node h in which the player had to play the actions associated x_1, x_2, \dots, x_i in that exact order to reach h . Multiple paths could fit to this description since the chance actions are not relevant here. Indeed, the value of y_i in this optimization program will represent the player’s contribution to the reach probability of such a node h from root h_0 if the player plays according to \mathbf{x} .

A.4 AMSGRAD

Reddi et al. (2018) proposed AMSGrad (AMS) as a fix to ADAM (Kingma & Ba, 2015), which may not converge in some stochastic convex optimization problems. AMS is described in Algorithm 5. The max operator, square root $\sqrt{\cdot}$, and division $/$ of vectors are to be interpreted element-wise. The projection operator $\Pi_{\Delta(m), \mathbf{v}} : \mathbb{R}^m \rightarrow \Delta(m)$ for a vector $\mathbf{v} \in \mathbb{R}_{\geq 0}^m$ is defined as

$$\mathbf{x} \mapsto \operatorname{argmin}_{\mathbf{y} \in \Delta(m)} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{v}} := \operatorname{argmin}_{\mathbf{y} \in \Delta(m)} \sqrt{\langle \mathbf{y} - \mathbf{x}, \mathbf{v}^T (\mathbf{y} - \mathbf{x}) \rangle}.$$

We implement this projection onto the simplex efficiently using the algorithm by Helgason et al. (1980).

A.5 DEPLOYMENT PHASE OF SIMULATION PROBLEMS

Figure 6 displays the subgame Γ' representing the deployment phase of the simulation problem we start to describe in Figure 2 (right).

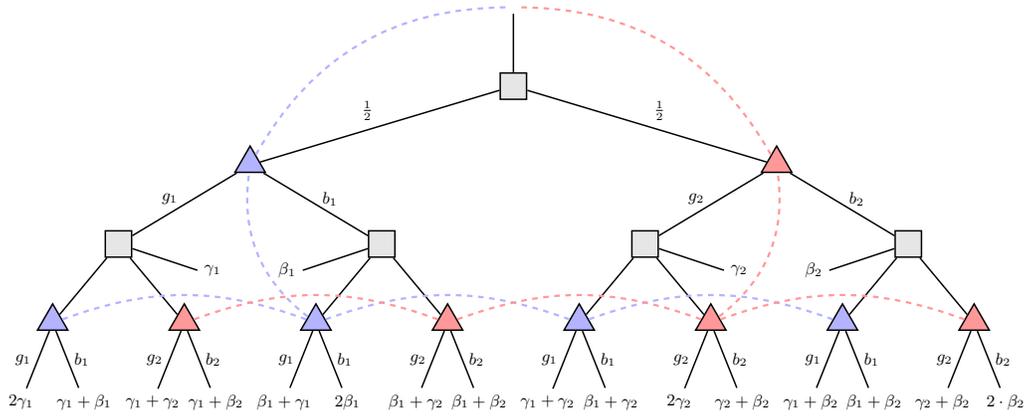
A.6 A PROOF OF PROPOSITION 4

We adopt the notation of Algorithm 1. For any $\mathbf{x}^* \in \mathcal{X}$, we have

$$\begin{aligned} \sum_{t=1}^T [U(\mathbf{x}^*) - U(\mathbf{x}^{(t)})] &\leq \sum_{t=1}^T \langle \nabla U(\mathbf{x}^{(t)}), \mathbf{x}^* - \mathbf{x}^{(t)} \rangle \\ &= \sum_{t=1}^T \sum_{i=1}^n \langle \mathbf{u}_i^{(t)}, \mathbf{x}_i^* - \mathbf{x}_i^{(t)} \rangle \leq O_T(\sqrt{T}). \end{aligned}$$

where we use, in turn: concavity of U , the definition of the utility vectors given to every information set, and the fact that RM is a regret minimizer (Hart & Mas-Colell, 2000). Dividing both sides by T and taking \mathbf{x}^* to be a global maximizer of U , it follows that there is always some iteration $1 \leq t \leq T$ for which $U(\mathbf{x}^*) - U(\mathbf{x}^{(t)}) \leq O_T(1/\sqrt{T})$.

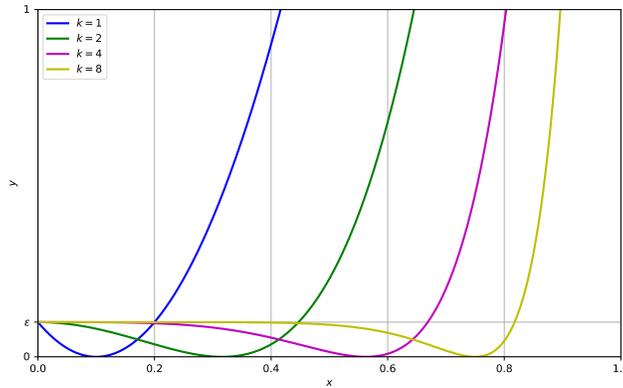
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039



1040
1041
1042
1043
1044

Figure 6: Deployment phase Γ' of the more complex simulation problem with two scenarios given in Figure 2 (right). In deployment, the agent acts at least once and up to two times in total. The “good” and “bad” actions yield different immediate payoffs in different scenarios, and they contribute additively to the total payoffs at terminal nodes.

1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058



1059
1060
1061
1062

Figure 7: A family of polynomial optimization problems over the unit interval on which the RM and GD families of algorithms perform arbitrarily poorly.

1063
1064

A.7 WHEN THE FIRST-ORDER OPTIMIZERS PERFORM POORLY

1065
1066
1067
1068
1069
1070
1071
1072

We find it quite surprising that the first-order optimizers we benchmark perform so well in terms of utility value in comparison to the global optimum found by `Gurobi`. Indeed, from a theoretical standpoint, we can give three examples in which the RM and/or GD families of algorithms converge to an arbitrarily bad value relative to the global optimum. To simplify explanations, we present all decision making under imperfect recall instances as maximization of a polynomial function over the 1-dimensional unit interval $[0, 1]$ instead. Theorem 1 describes how to efficiently construct the decision making under imperfect recall instance from that.

1073
1074
1075
1076
1077
1078
1079

Example 1: All converge to a bad value Consider the (ϵ, k) -parametrized function $f_{\epsilon, k}(x) = \frac{1}{\epsilon}(x^k - \epsilon)^2$ for $\epsilon > 0$ and $k \in \mathbb{N}$ over the unit interval $[0, 1]$. The function f is plotted in Figure 7 for $\epsilon = 0.1$ and multiple values for k . In all cases, $f_{\epsilon, k}(x) \geq 0$, $f_{\epsilon, k}(0) = \epsilon$, and $f_{\epsilon, k}(1) > 1$ if we additionally restrict $\epsilon < \frac{3-\sqrt{5}}{2} \approx 0.382$. If an algorithm therefore converges to $x^* = 0$, we have found a family of instances for which the algorithm has achieved no more than $\text{MIN} + \epsilon \cdot (\text{MAX} - \text{MIN})$ in value, where MAX and MIN represent the max and min values f on $[0, 1]$ (or, respectively, the utility function on the 1-simplex). For our first-order methods, note that $f_{\epsilon, k}$ is strictly decreasing

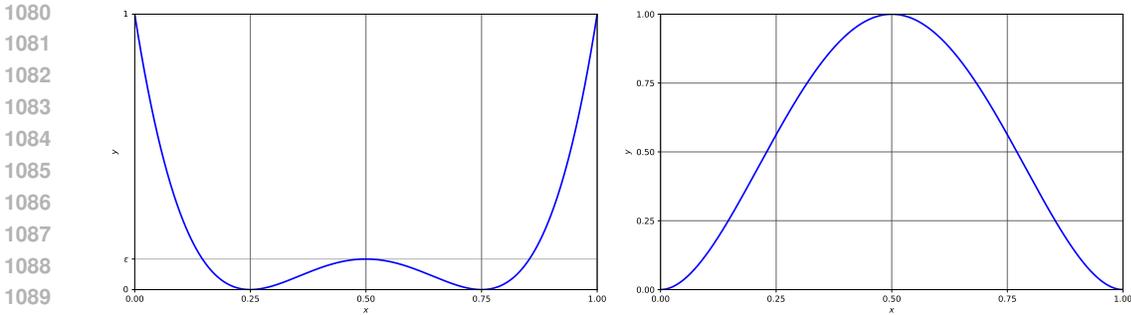


Figure 8: Examples in which particular first-order optimizers are likely to perform arbitrarily bad in value while others converge to the global optimum (if not initialized at $x^* = 1/2$ exactly). On the left, the gradient descent methods initialized closer to the middle will converge to the local optimum at $x^* = 1/2$ while the regret matching methods converge to the global optima. On the contrary, regret matching methods on the right will converge to the global minimum (which is also first-order maximal) while the gradient descent methods converge to the global optimum.

in the interval $J = [0, \epsilon^{1/k}]$. If the RM and GD families of algorithms are therefore initialized to start in J , they will converge to 0. Assuming we draw the initial point uniformly random from $[0, 1]$, this situation occurs with probability $\epsilon^{1/k}$. Therefore, we can first set the desired poor-performance parameter ϵ , and then $k = k(\epsilon)$ to meet the desired probability confidence $\epsilon^{1/k}$, to obtain arbitrarily bad performance of the RM and GD families of algorithm with arbitrarily high probability in some instance.

Example 2: gradient descent methods converge to a bad value Consider the function $f(x) = (\frac{1}{4} \cdot \frac{3}{4})^2 \cdot (x - \frac{1}{4})^2 \cdot (x - \frac{3}{4})^2$ over the unit interval $[0, 1]$, as plotted in Figure 8 (left). Then GD, OGD, and AMS for any step size ν will converge to the local maximum $x^* = 1/2$ if initialized sufficiently close to it. We omit the extension of this singular example to a parametrized family of example—for which we would work with exponents of f similar to Example 1 above—in order to increase the neighborhood of attraction around $x^* = 1/2$ and increase the difference $\frac{f(1/2)}{f(0)}$ of local optimum to global optimum. Independent of such extensions, the regret matching methods will always take such a large step at the first iteration that it reaches the global optimum at $x \in \{0, 1\}$ at the second iteration and stays there thereafter; that is, as long as the starting point is not exactly $x^* = 1/2$ or exactly one of the stationary points with value 0 (for f , that is $\{\frac{1}{4}, \frac{3}{4}\}$).

Example 3: Regret matching methods converge to a bad value We consider the fact that regret matching methods start off with very large steps in the direction of the gradients a blessing rather than a curse for the convergence speed of the first-order optimizer. However, we can also exploit that property when it comes to performance in terms of value achieved. Consider the function $f(x) = 16 \cdot x^2(1 - x)^2$ over the unit interval $[0, 1]$, as plotted in Figure 8 (right). Then as long as the regret matching methods do not start at the global optimum $x^* = 1/2$ exactly, they will grossly overshoot their first step in the direction of the gradient, and reach $x^{(2)} \in \{0, 1\}$ in the second iteration. Despite yielding the globally minimum value of 0, the regret matching methods will be stuck there thereafter since those two points are also stationary points. The gradient descent methods for any step size ν , on the other hand, will converge to the global optimum at $x^* = 1/2$ as long as they start sufficiently close it. Again, we omit the extension of this example to a parametrized class of examples that show how one can increase the neighborhood of attraction around $x^* = 1/2$.

A.8 ADDITIONAL EXPERIMENTAL DETAILS AND RESULTS

All experiments were run on a 64-core AMD Opteron 6272 processor. Each run was allocated one thread with a maximum of 16GBs of RAM. The commercial solver Gurobi requires a license to run on decision problems of nontrivial size. The result table of the experiments for the full set of benchmark decision problems is given in Table 3, which now also includes PRM experiments. We display “—” in the time column of Gurobi if it does not converge to the global optimum (up to a

Problem	Gurobi		GD		OGD		AMS		RM		RM*		PRM		PRM*								
	value	time	value	gap	value	gap	value	gap	value	gap	value	gap	value	gap	value	gap							
Det-86	18.00	0.22s	18.00	0.01s	18.00	0.01s	18.00	0.06s	18.00	0.00s	18.00	0.00s	18.00	0.00s	18.00	0.00s							
Det-105	12.00	1.86s	12.00	0.01s	12.00	0.01s	12.00	0.05s	12.00	0.00s	12.00	0.00s	12.00	0.00s	12.00	0.00s							
Det-1k	13.00	1m 24s	13.00	0.13s	13.00	0.07s	13.00	1.04s	13.00	0.32s	13.00	0.36s	13.00	0.38s	13.00	0.41s							
Det-1.8k	22.00	2m 40s	22.00	0.06s	22.00	0.07s	22.00	0.71s	22.00	0.03s	22.00	0.03s	22.00	0.03s	22.00	0.03s							
Det-2.0k	17.50	1m 42s	17.50	0.03s	17.50	0.05s	17.50	0.20s	17.50	0.03s	17.50	0.03s	17.50	0.03s	17.50	0.03s							
Det-8k	16.67	—	16.62	13.05s	16.62	2m 36s	16.67	8m 1s	16.67	—	3e-05	16.67	2m 39s	16.67	—	0.001	16.67	—	0.007				
Det-10.6k	12.84	—	12.70	24.87s	12.70	5m 0s	12.84	4m 44s	12.84	6.41s	12.84	6.74s	12.84	15.53s	12.84	14.48s	—	—	—				
Det-10.7k	20.20	16m 36s	20.20	0.21s	20.20	0.23s	20.20	2.40s	20.20	0.73s	20.20	0.77s	20.20	1.06s	20.20	1.13s	—	—	—				
Det-86k	14.89	—	14.84	—	0.004	10.00	—	11.4	14.89	—	7e-05	14.89	2m 7s	14.89	2m 4s	14.89	7m 54s	14.89	5m 50s				
Det-130k	15.53	—	15.37	—	8e-06	15.40	45m 59s	—	15.53	57m 57s	—	15.53	5m 38s	15.53	2m 3s	15.53	—	0.0001	15.53	—	8e-05		
Det-139k	18.89	—	18.76	15m 24s	—	18.76	16m 48s	—	18.89	31m 26s	—	18.89	1m 47s	18.89	1m 50s	18.89	3m 29s	18.89	3m 10s	—	—		
Det-718k	—	—	12.76	—	0.0005	12.69	—	0.005	12.84	2h 9m	—	12.84	30m 55s	12.84	31m 21s	12.84	—	0.001	12.84	—	0.0008		
Det-1.002m	—	—	13.93	—	0.0003	13.90	—	0.005	13.96	—	1e-05	13.96	15m 36s	13.96	17m 17s	13.96	40m 10s	13.96	35m 35s	—	—		
Det-1.008m	—	—	12.64	—	0.0006	12.54	—	0.008	12.75	—	3e-05	12.75	33m 31s	12.75	54m 51s	12.75	31m 28s	—	—	—	—		
Det-2.1m	—	—	26.00	—	1e-05	25.96	—	0.02	26.15	—	0.002	26.15	—	0.003	26.15	3h 25m	26.15	—	0.006	26.15	—	0.005	
Det-2.2m	—	—	16.20	—	0.002	15.93	—	0.02	16.36	—	0.0002	16.36	2h 22m	—	16.36	3h 13m	16.36	—	2e-06	16.36	—	5e-06	
Det-3.8m	—	—	15.66	—	0.003	15.14	—	0.03	15.78	—	0.0002	15.80	—	2e-06	15.80	—	5e-05	15.80	—	0.002	15.80	—	0.0003
Det-4.0m	—	—	18.17	—	0.005	17.72	—	0.03	18.33	—	0.0005	18.34	—	2e-05	18.34	2h 55m	18.34	—	0.005	18.34	—	0.005	
Det-4.1m	—	—	17.88	—	0.003	17.47	—	0.03	18.05	—	0.0004	18.06	—	4e-05	18.06	—	2e-05	18.06	—	0.003	18.06	—	0.0007
Det-4.2m	—	—	19.98	—	0.003	20.07	—	0.003	20.15	—	0.0004	20.15	—	2e-05	20.15	—	0.01	20.15	—	0.02	—	—	
Det-9m	—	—	23.16	—	0.004	22.71	—	0.02	23.45	—	0.004	23.45	—	0.0001	23.45	—	0.0003	23.45	—	0.0003	23.45	—	0.0004
Det-10m	—	—	24.64	—	0.002	24.61	—	0.003	24.76	—	0.009	24.76	—	0.002	24.76	—	0.0004	24.76	—	0.01	24.76	—	0.0008
Det-18m	—	—	26.38	—	0.006	25.81	—	0.05	26.71	—	0.004	26.71	—	0.004	26.71	—	0.001	26.71	—	0.04	26.71	—	0.04
Rand-7k	0.53	25m 18s	0.49	4.88s	—	0.49	5.14s	—	0.50	2.69s	—	0.50	0.38s	—	0.50	0.27s	—	0.50	0.26s	—	0.50	0.34s	—
Rand-11.9k	1.00	1h 16m	0.97	0.93s	—	0.97	0.90s	—	0.99	1.38s	—	0.95	0.26s	—	0.95	0.29s	—	0.95	0.19s	—	0.95	0.23s	—
Rand-12.2k	1.00	1h 52m	0.93	3.33s	—	0.92	2.73s	—	0.92	2.35s	—	0.93	0.36s	—	0.93	0.41s	—	0.94	0.36s	—	0.94	0.41s	—
Rand-24k	0.72	—	0.66	7m 0s	—	0.66	7m 46s	—	0.66	4m 4s	—	0.66	26.55s	—	0.66	1m 3s	—	0.66	1m 54s	—	0.66	5m 5s	—
Rand-35k	1.00	—	0.95	3.85s	—	0.95	3.76s	—	0.95	3.90s	—	0.92	0.99s	—	0.92	1.18s	—	0.92	0.92s	—	0.94	1.68s	—
Rand-42k	0.69	—	0.55	—	0.01	0.55	—	0.01	0.64	—	0.0006	0.65	—	2e-06	0.65	5m 56s	—	0.65	—	5e-06	0.65	3m 19s	—
Rand-165k	0.37	—	0.96	19.77s	—	0.97	18.48s	—	0.99	28.90s	—	0.96	4.33s	—	0.97	4.95s	—	0.96	5.24s	—	0.90	4.02s	—
Rand-179k	0.38	—	0.88	—	0.0003	0.88	—	1e-06	0.96	1m 7s	—	0.94	5.97s	—	0.93	10.27s	—	0.93	6.66s	—	0.91	7.31s	—
Rand-198k	0.40	—	0.96	25.37s	—	0.95	22.61s	—	0.97	39.73s	—	0.96	8.10s	—	0.96	7.41s	—	0.95	5.22s	—	0.96	6.31s	—
Rand-1.2m	—	—	0.93	2m 46s	—	0.93	2m 28s	—	0.97	2m 17s	—	0.96	35.86s	—	0.97	36.07s	—	0.96	31.74s	—	0.96	31.09s	—
Rand-1.3m	—	—	0.96	4m 0s	—	0.96	3m 17s	—	1.00	8m 33s	—	0.96	2m 26s	—	0.96	54.77s	—	0.98	1m 33s	—	0.93	38.99s	—
Rand-2m	—	—	0.92	3m 53s	—	0.93	3m 44s	—	0.97	2m 37s	—	0.94	59.29s	—	0.93	1m 21s	—	0.96	2m 1s	—	0.96	58.84s	—
Rand-4m	—	—	0.94	13m 1s	—	0.94	15m 39s	—	0.96	29m 2s	—	0.93	3m 12s	—	0.92	4m 26s	—	0.92	8m 41s	—	0.93	4m 16s	—
Rand-6m	—	—	0.97	17m 34s	—	0.97	15m 40s	—	0.99	14m 23s	—	0.98	2m 30s	—	0.98	2m 9s	—	0.98	2m 50s	—	0.98	2m 10s	—
Rand-7m	—	—	0.97	22m 27s	—	0.98	25m 7s	—	0.99	11m 45s	—	0.94	2m 13s	—	0.93	2m 52s	—	0.96	2m 55s	—	0.97	3m 47s	—
Rand-13m	—	—	0.59	—	0.003	0.58	—	0.003	0.65	1h 40m	—	0.63	19m 11s	—	0.64	17m 31s	—	0.64	20m 39s	—	0.65	36m 42s	—
Rand-18m	—	—	0.97	2h 33m	—	0.97	3h 0m	—	0.99	1h 29m	—	0.95	29m 45s	—	0.97	24m 0s	—	0.96	13m 8s	—	0.97	14m 31s	—
Rand-23m	—	—	0.94	3h 37m	—	0.93	—	0.0007	0.98	3h 20m	—	0.98	23m 10s	—	0.96	23m 5s	—	0.98	16m 48s	—	0.95	18m 2s	—
Sim-245	4.41	0.18s	4.41	0.00s	4.41	0.00s	4.41	0.01s	4.41	0.01s	—	4.41	0.00s	4.41	0.00s	4.41	0.00s	4.41	0.00s	4.41	0.00s	4.41	0.00s
Sim-438	7.21	0.41s	7.21	0.00s	7.21	0.00s	7.21	0.01s	7.21	0.01s	—	7.21	0.00s	7.21	0.00s	7.21	0.00s	7.21	0.00s	7.21	0.00s	7.21	0.00s
Sim-759	3.89	2.97s	3.89	0.01s	3.89	0.01s	3.89	0.02s	3.89	0.01s	—	3.89	0.01s	3.89	0.01s	3.89	0.01s	3.89	0.01s	3.89	0.01s	3.89	0.01s
Sim-3k	6.25	1m 1s	6.25	0.32s	6.25	1.03s	6.25	5.54s	6.25	0.26s	—	6.25	0.28s	6.25	0.52s	6.25	0.52s	6.25	0.48s	6.25	0.48s	6.25	0.48s
Sim-7k	8.58	1m 36s	8.58	0.05s	8.58	0.05s	8.58	0.12s	8.58	0.05s	—	8.58	0.05s	8.58	0.05s	8.58	0.05s	8.58	0.05s	8.58	0.05s	8.58	0.05s
Sim-13k	10.38	4m 21s	10.38	0.69s	10.38	8.54s	10.38	14.37s	10.38	1.03s	—	10.38	1.01s	10.38	4.75s	10.38	4.75s	10.38	3.97s	10.38	3.97s	10.38	3.97s
Sim-34k	10.44	1h 42m	10.44	4.89s	10.44	6.74s	10.44	1m 9s	10.44	2.52s	—	10.44	2.81s	10.44	5.03s	10.44	5.03s	10.44	5.01s	10.44	5.01s	10.44	5.01s
Sim-66k	6.94	1h 31m	6.94	5.63s	6.94	8.70s	6.94	1m 9s	6.94	5.51s	—	6.94	3.94s	6.94	17.32s	6.94	17.32s	6.94	15.01s	6.94	15.01s	6.94	15.01s
Sim-105k	4.40	—	4.40	18.60s	4.40	1m 0s	4.40	2m 35s	4.40	2m 41s	—	4.40	55.90s	4.40	15m 18s	4.40	15m 18s	4.40	10m 51s	4.40	10m 51s	4.40	10m 51s
Sim-125k	14.47	—	14.48	12.70s	14.48	19.76s	14.48	4m 46s	14.48	11.70s	—	14.48	12.15s	14.48	18.83s	14.48	18.83s	14.48	19.68s	14.48	19.68s	14.48	19.68s
Sim-226k	8.57	—	9.70	2.16s	9.70	4.28s	9.70																