

COGNITIVE TRAINING FOR LANGUAGE MODELS: TOWARDS GENERAL CAPABILITIES VIA CROSS-ENTROPY GAMES

CLÉMENT HONGLER^{1,2} FRANCK GABRIEL^{3,1} VALENTIN HARTMANN¹ ARTHUR RENARD¹ ANDREW EMIL¹
¹XENT LABS ²EPFL ³UNIVERSITÉ LYON 1

ABSTRACT. Defining a constructive process to build general capabilities for language models in an automatic manner is considered an open problem in artificial intelligence. Towards this, we consider the problem of building a curriculum of tasks that grows a model via *relevant skill discovery*.

We provide a concrete framework for this task, using a family of tasks called Cross-Entropy Games, which we postulate is universal in a suitable sense. We show that if it is possible to grow the curriculum for relevant skill discovery by iterating a greedy optimization algorithm, then, under natural assumptions, there is essentially only one *meta-objective* possible (up to a few hyper-parameters). We call the resulting process *cognitive training*.

We postulate that, given sufficiently capable language models as players and meta-samplers, cognitive training provides a principled way to relevant skill discovery; and hence to the extent general capabilities are achievable via greedy curriculum learning, cognitive training would be a solution.

1. ARTIFICIAL GENERAL INTELLIGENCE AND COGNITIVE TRAINING

The last decades saw spectacular progress on several fronts for AI. In particular we saw:

- Models that could generalize beyond their training dataset.
- Reinforcement learning that could reach super-human performance on specific tasks.
- Pre-trained Large Language Models (LLMs) that could learn to perform certain tasks with in-context learning, and numerous more specific tasks, given a suitable training environment.

In spite of these achievements, there is a growing consensus that a number of ideas are lacking to go towards “true Artificial General Intelligence (AGI)”. In this paper, we start with a short summary of what we expect from AGI starting from seminal works [Tur50, LeHu07] and propose to study an (apparently simpler) problem, which is that of (open-ended) relevant skill discovery.

To study this problem, we propose a training framework (based on an elaboration of the framework introduced in [HoEm25]) which we call *cognitive training*: the idea is to grow a curriculum of games in a specific space of tasks, called *Cross-Entropy (Xent) Games*, endowed with a suitable notion of *transfer value*: how learning one game teaches one to play another. We postulate that the learning achieved by any curriculum of tasks can be approximated by a curriculum of Xent Games.

We consider the problem of growing a curriculum of Xent Games in a greedy manner, by optimizing a meta-objective \mathcal{O} . We show that (perhaps surprisingly), given a few natural assumptions, a consistent meta-objective must take a very constrained form: from these, we derive an explicit formula for \mathcal{O} , that balances sparsity (inverse code length), quality (internal consistency), diversity (novelty), and external relevance (benchmark performance).

This leads us to *cognitive training*. The idea is to start with one (or more) base auto-regressive model \mathcal{M} , find a suitable space of Xent Games \mathcal{G} . At every step of the curriculum, a meta-sampler \mathcal{M}_S then generates games on \mathcal{G} to optimize \mathcal{O} , to grow a curriculum that brings relevant skill discovery to \mathcal{M} .

Our reasoning thus leads (from a number of assumptions) to the following idea: *if it is possible to greedily generate a curriculum of tasks bringing general capabilities to a model, then it is possible to replicate this with a greedily built curriculum of Xent Games, and then the meta-objective formula must take a very specific form.*

1.1. Definition of Artificial General Intelligence. The first notable attempt at a measure of artificial intelligence is probably Turing’s imitation game [Tur50], leading to the famous Turing test. Across decades, various definitions of intelligence in the context of AI were provided, leading to a notable and influential synthesis in [LeHu07]:

Intelligence measures an agent’s ability to achieve goals in a wide range of environments.

S. Legg and M. Hutter

A framework built at the intersection of Occam’s Razor (in the form of Algorithmic Probability) and of Reinforcement Learning [LeHu06] highlighted notably a theoretically well-posed, but uncomputable, approach to the so-called Universal AGI problem: the AIXI paradigm [Hut05], blending a universal prior on world models (based on Kolmogorov’s complexity) and objective maximization.

Arguably, the part that remains a little vague is associated with the “wide range” in the above definition: is there a good measure of generality? Are humans general according to any such definition? Is there a process through which this generality is discovered?

Despite the impossibility to apply this framework directly (due to its uncomputability) and the lack of answers about generality, AIXI gives a direction for what we would expect theoretically from an AI going towards AGI: models should interact with an environment, perform tasks within them, and be rewarded. Two decades of progress in the field of Reinforcement Learning (RL) followed. In the context of language models, this was discussed for quite some time [MJB15]. However, as discussed in the next subsection, the progress has only recently accelerated thanks to the rise in prominence of pre-trained large language models.

1.2. Cognitive Training: Goals. The current problem with language models is not that they cannot learn tasks given a suitable RL environment, but that they have a hard time performing (or quickly learning to perform) tasks for which they have not been trained already. Given a new task, constructing a specifically relevant RL environment and post-training a model on it can be very difficult or impossible.

The question that we aim to address can be understood as building extended pre-training using *games* (tasks/environments with an emphasis on certain types of rewards):

Problem 1. Can we find a curriculum of compact text games $(G_k)_{k \geq 0}$ such that a model \mathcal{M} trained on it is maximally ready for new/unseen tasks given some data, architecture, and compute constraints?

Mathematically, the above problem seems to assume a prior on new tasks, and is hence somewhat ill-posed (we have no way to set this prior). Like for supervised learning (where e.g. computer vision problems could be solved without a clear prior on the set of images corresponding to a given label), it is reasonable to expect that one can achieve good performance with the right methods, without expliciting (or sampling from) such a prior. By Occam’s razor, if we find a curriculum of games (at finite k) that yields improved performance on an external metric E (e.g. an aggregation of benchmarks), then it is reasonable to expect that this curriculum will be *generally* valuable: a model trained on such games is likely (in an informal sense) to generalize well to tasks beyond the specific tasks of E .

Pre-training is in fact the most emblematic example of a compact training game: predicting the next token, a most universal and abstract task, leads us quite far in terms of general capabilities. Given a dataset, this is, however, a game that can only be played once (and with diminishing returns) leading to the question: what are, besides pre-training, the most useful training games?

Interestingly, we can perform a substantial reduction to the above question by asking the following, which seems perhaps a simpler problem a priori:

Problem 2. Can we find a curriculum of compact games $(G_k)_{k \geq 0}$ such that an LLM \mathcal{M} will keep discovering relevant new skills throughout training, while maintaining existing capabilities?

The following seems intuitive, as any solution to Problem 1 would (at least eventually) involve the discovery of new relevant skills (which are needed to deal with new tasks):

Claim 3. Any convincing solution to Problem 1 should also yield a solution to Problem 2.

According to Claim 3, finding a curriculum as described in Problem 2 is a prerequisite for training a model that is maximally prepared for new tasks. Deriving a process for finding such a curriculum is the goal of this work.

To make the above problem concrete, we propose a space of games, called *Cross-Entropy Games (Xent Games)*, endowed with an appropriate structure of transfer value; separately, it is postulated that any curriculum of tasks can be approximated by a curriculum of Xent Games (Section 3.3); this leads to a specific framework for the above definition (Definition 19).

We then focus on *greedy curriculum building* (i.e. where each game is picked by optimizing a function of the previous games). This results in the central (and perhaps most surprising) claim of this paper (Sections 4.2–4.3):

Claim 4. If we can solve Problem 2 via a greedy method involving at every step the optimization of meta-objective \mathcal{O} , then from reasonable principles an explicit (and unique) formula for \mathcal{O} can be found.

This leads us to the concept of *cognitive training*: the process of performing the meta-optimization for \mathcal{O} . In the next subsection, we first lay out a framework to optimize cognitive training for LLMs, detailed in the rest of the paper.

1.3. Cognitive Training: Implementation. The question raised in Problem 1 above suggests formulating a meta-optimization problem: that of building a curriculum of games that turn out to be maximally useful to a (language) model \mathcal{M} learning to play them (for a good definition of “usefulness”), on which one trains a model sequentially.

To implement this program, we need two ingredients:

- a space of games \mathcal{G} on which the model \mathcal{M} will be trained,
- a meta-objective \mathcal{O} defined on games \mathcal{G} , used to select which games to add to the curriculum.

Given these two ingredients, \mathcal{O} is optimized at each step by a meta-sampler \mathcal{M}_S , which outputs games (written in a specific language). The next game added to the curriculum to train model \mathcal{M} is the one that maximizes the meta-objective \mathcal{O} . In the rest of this subsection, we briefly discuss principles for the design of \mathcal{G} and \mathcal{O} ; details are presented in Sections 3 and 4, respectively.

1.3.1. Game Space Design Elements. The desiderata for the space of games \mathcal{G} aimed at developing the cognitive capabilities of a model \mathcal{M} are numerous:

- The games should be amenable to fast training.
- It should be easy to generate a wide diversity of games in \mathcal{G} with fairly concise code.
- The space should allow for robust transfer exploration: if a certain curriculum (based on certain games not necessarily in \mathcal{G}) teaches \mathcal{M} some skills, the same can be achieved by a curriculum of games in \mathcal{G} .

In Section 2 below, we argue that a fertile ground for building the games in \mathcal{G} is to rely on the *implicit knowledge* of LLMs; this leads to the proposal to use Cross-Entropy Games (Xent Games) as a means from which to build \mathcal{G} .

1.3.2. Meta-Game Design Elements. Our main contribution is the principled derivation in Section 4 of a meta-algorithm, called *cognitive training*, to build a useful curriculum of games $(G_k)_{k \geq 0}$ in \mathcal{G} to develop the cognitive capabilities of a model \mathcal{M} , building a sequence of models where \mathcal{M}_{k+1} is obtained by training \mathcal{M}_k on G_k . The *cognitive training* algorithm is a greedy algorithm based on a *meta-objective* \mathcal{O} . Central to the definition of \mathcal{O} is the concept of *transfer value* (Section 3.4.2), which estimates the relevance of learning a game to play another.

At each step $k \geq 1$, the choice of a new game H should balance:

- quality/relevance q : how much H improves performance on the old games $G_{<k} := G_1, \dots, G_{k-1}$;
- diversity/novelty d : how much H brings skills that are not yet captured by the games $G_{<k}$;
- external benchmark performance b : how much H improves performance on external metrics;
- description length l : the raw code length of H .

In Section 4, we derive explicit principled definitions of each of the previous quantities, and derive a formulation for the meta-objective \mathcal{O} , namely

$$\mathcal{O} = \frac{qd + bp}{l},$$

where p is a *pressure* factor, balancing the relative importance of internal and external performance. Note that despite involving a priori a number of terms that scale linearly in k (and similarly for memory), the meta-objective \mathcal{O} satisfies some good scalability properties (see Section 4.3.6).

1.3.3. Scope and Contributions. The key contributions of this paper are the following:

- A treatment of Problem 2, i.e. relevant skill discovery, with the idea that if general capabilities are achieved by curriculum learning, they automatically imply relevant skill discovery.
- A formal definition, building upon [HoEm25], of a space of tasks called *(Streamlined) Cross-Entropy (Xent) Games*, endowed with a transfer value structure.
- A hypothesis that curricula made of (Streamlined) Xent Games can approximate the learning of any curriculum; and hence that relevant skill discovery can be understood in terms of Xent Games and their transfer value structure.
- A formalization of greedy curriculum learning for Xent Games in terms of a meta-objective optimization.
- A derivation of an explicit formula for any meta-objective satisfying suitable consistency assumptions.

1.3.4. Structure of the Paper.

- In Section 2, we outline the key idea of implicit knowledge for an LLM: this suggests that much more can be done with an LLM than simple sampling.
- In Section 3, we present the space of games upon which cognitive training is based: the (Streamlined) Xent Games, derived from the implicit knowledge of LLMs.
- In Section 4, we derive a natural framework for xent-game-based cognitive training for LLMs.
- In the same Section 4, we discuss several challenges and open questions associated with cognitive training.

2. IMPLICIT KNOWLEDGE

If the models that we train are supposed to be useful for general tasks relevant to some world (e.g. human society), the games that they are trained on need to have a connection to that world. For language models, we posit that this connection is provided by the *implicit knowledge* of Large Language Models (LLMs) that have been trained on vast corpora of text data.

LLMs are generative models: given a context c , they define a conditional probability measure $\mathbb{P}_{\mathcal{M}}(\cdot | c)$ on sequences of tokens. For a given model \mathcal{M} , the likelihood is directly accessible: for any token sequence $x = x_1, \dots, x_T$, the quantity $\log \mathbb{P}_{\mathcal{M}}(x | c) = \sum_{t=1}^T \log \mathbb{P}_{\mathcal{M}}(x_t | c, x_{<t})$ is readily available from \mathcal{M} . As a consequence, the information we can extract from \mathcal{M} depends on how we use it: either as a generator, or as the basis for an algorithm relying on $\mathbb{P}_{\mathcal{M}}$ for e.g. perplexity comparison, search, and optimization in the space of token sequences.

In this section, we discuss the difference between:

- Explicit knowledge $\mathcal{E}_{\mathcal{M}}$: the information that can be quickly and reliably extracted from \mathcal{M} by sampling answers from well-chosen prompts.
- Implicit knowledge $\mathcal{I}_{\mathcal{M}}$: the information that can be extracted by arbitrary algorithms with access to $\mathbb{P}_{\mathcal{M}}$.

Since generation is itself a (randomized) algorithm relying on $\mathbb{P}_{\mathcal{M}}$, we have $\mathcal{E}_{\mathcal{M}} \subset \mathcal{I}_{\mathcal{M}}$. In the following, we will see that there is a gap between the two kinds of knowledge: focusing on continuations allows one to explore only a tiny subset of observables (functions) of the measure $\mathbb{P}_{\mathcal{M}}$, while direct access to the model’s measure allows for different procedures (search, contrastive experiments, ...) that further reveal the cognitive capabilities of the LLM. In Section 3, we exploit the implicit/explicit gap to define the Xent Games, designed to use implicit capabilities to improve the explicit knowledge of LLMs.

2.1. Explicit Knowledge. The *explicit knowledge* of an LLM is defined as follows: the set of questions and tasks for which, given a fixed computation budget, a model \mathcal{M} , typically fine-tuned for instruction following, can provide a correct solution with sufficiently high probability by generating a continuation. This knowledge is essentially what a standard benchmark (e.g. one based on the samples generated by the LLM) would measure.

The notion of explicit knowledge has some important subtleties:

- (1) The explicit knowledge for raw pre-trained models (not instruction-tuned) is ill-defined or poorly representative of the underlying knowledge of the LLM: for such models, the main objective is not to answer or to follow the output specifications.¹
- (2) Explicit knowledge depends on the decoding policy: it measures both the model distribution and an external choice of how we generate from it (greedy/sampling, temperature, ...).
- (3) The answers of an LLM can be stochastic, so we may allow for the sampling of a few answers before an algorithm produces the final output from these samples.

Finally, fine-tuning can modify part of the explicit knowledge: for instance, after alignment training, a model may refuse to answer a question its base version could have answered. This motivates the following perspective for understanding LLMs’ capabilities: instead of asking the model what it knows, we can ask what can be extracted from the model probabilities.

2.2. Implicit Knowledge. We define the *implicit knowledge* as follows: the set of questions/tasks that can be answered algorithmically (given a fixed computation budget) given access to the model measure $\mathbb{P}_{\mathcal{M}}$. More specifically, we assume that we have access to the log-likelihoods of continuations $\log \mathbb{P}_{\mathcal{M}}(x | c)$. The notion of implicit knowledge is obviously appealing as it is the set of information that an LLM already knows (in one form or another) from its pre-training.

¹With well-chosen prompts, one can steer LLMs towards the desired behavior, but this will not replace a good fine-tuning.

The implicit knowledge differs from the explicit knowledge in at least three important ways. First, it is well defined for raw pre-trained models, as it does not rely on a question-answering setting. Second, it does not depend on the decoding policy as it depends solely on $\mathbb{P}_{\mathcal{M}}$. Third, it is, by its nature, related to a deterministic computational problem, rather than the consequence of a stochastic generator.

The simplest but most fundamental implicit knowledge is the ability to provide the log-likelihood of a sentence or continuation. For any $x = x_1 \dots x_T$,

$$\log \mathbb{P}_{\mathcal{M}}(x_1 \dots x_T | c) = \sum_{t=1}^T \log \mathbb{P}_{\mathcal{M}}(x_t | c, x_{<t}),$$

can be extracted from the model’s measure, and is therefore part of the implicit knowledge. Yet, this quantity is not part of the explicit model: a model generally cannot explicitly report its correct value.

Of course, by sampling more and more continuations, we should theoretically be able to recover $\log \mathbb{P}_{\mathcal{M}}(x | c)$. Yet, the number of samples required is prohibitive. Indeed, for a fixed continuation x , observing x even once requires an order of $1/\mathbb{P}_{\mathcal{M}}(x|c)$ samples. Since $\mathbb{P}_{\mathcal{M}}(x | c)$ typically decreases exponentially with the length of x , it is impossible to recover the log-likelihood from samples within a reasonable computing time. This simple argument shows that having direct access to the log-likelihood is qualitatively different from having access only to the generation process.

2.3. Examples of Implicit Knowledge of Pre-Trained Models. Direct access to the log-likelihood enables procedures that are out of reach for explicit generation. As explained above, this is important for raw pre-trained models: even when a model does not reliably follow a question/answer format, its conditional measure can still be evaluated through log-likelihoods. In this section, we provide a non-exhaustive family of procedures for extracting implicit knowledge.

One way to do so is to consider the difference between log-likelihoods in which some candidate sequences x_1 or x_2 are inserted into possibly different contexts c_1 and c_2 . This quantity measures how much more the model supports x_1 in the hypothetical world specified by c_1 than x_2 in the hypothetical world specified by c_2 . By choosing contexts appropriately and determining how to insert the candidate sequences into them, this likelihood difference can be used among other things to extract beliefs from a model, quantify the usefulness of some information, or define some contrastive objectives. In practice, contexts are built from a base context to which pre-prompts or sandwich prompts are added before and after it; we omit these details here, but they matter in concrete implementations.

2.3.1. Likelihood difference for a fixed statement. We can compare the likelihood of the same statement in two different contexts.

- **Truth-false difference.** For a statement s , we can consider

$$\begin{aligned} \Delta_{\text{T/F}}(s | c) &= \log \mathbb{P}_{\mathcal{M}}(s | c, \text{”The following statement is true : ”}) \\ &\quad - \log \mathbb{P}_{\mathcal{M}}(s | c, \text{”The following statement is false : ”}) \end{aligned}$$

A simple classifier predicts “true” whenever $\Delta_{\text{T/F}}(s | c) > 0$. To reduce dependence on the exact prompt, we can compute the average of $\Delta_{\text{T/F}}$ over paraphrases of the prefix and/or of the statement s . Closely related, in [KCAHK22], to study the calibration of models, the authors present a proposed answer to a model and ask it whether it is true or false, then read off the log-likelihood of the two options.

- **Multiple-choice selection for explanations.** More generally, given candidate explanations c_1, \dots, c_n , we can select $\arg \max_{i \in \{1, \dots, n\}} \log \mathbb{P}_{\mathcal{M}}(s | c_i)$. This allows one to select the most plausible explanation for s .

2.3.2. Counterfactual thinking and surprise reduction. A variant of the multiple-choice selection of explanations is to quantify the value of some information. Let q be a problem, x a reference solution, and h an additional piece of information. The difference

$$\Delta_{\text{info}}(h; q \rightarrow x) = \log \mathbb{P}_{\mathcal{M}}(x | q, h) - \log \mathbb{P}_{\mathcal{M}}(x | q)$$

measures how much h reduces the model’s surprise about x . This makes it possible to quantify the usefulness of a piece of information without requiring the model to explicitly explain why (which does not provide a reliable, quantitative value).

2.3.3. *Contrastive objective across models or checkpoints.* Another variant of log-likelihood comparison is to consider two different models or checkpoints. Given two models \mathcal{M}_1 and \mathcal{M}_2 , the quantity

$$\Delta_{\text{ctr}}(x | c) = \log \mathbb{P}_{\mathcal{M}_1}(x | c) - \log \mathbb{P}_{\mathcal{M}_2}(x | c)$$

quantifies whether x is more plausible for \mathcal{M}_1 than for \mathcal{M}_2 . When \mathcal{M}_1 is stronger, for example if it is a larger model, and \mathcal{M}_2 is weaker, this quantifies the continuations that are more interesting in the sense that a clever model can think about them, but a simpler model cannot. This is closely related to contrastive sampling methods (see e.g. [LHFLL22]).

2.3.4. *Verification vs construction.* Many tasks have solutions that are easier to verify than to construct.² A well-trained model can assign a higher likelihood to correct solutions than to incorrect ones, while still failing to construct a correct solution by direct generation.

Access to the log-likelihood transforms the verification task into a continuous signal: rather than a discrete “correct/incorrect” reward, $\log \mathbb{P}_{\mathcal{M}}$ can be used to construct a continuous objective that can guide search and optimization.

In addition, it is worth noting that some continuity in reward signals can be provided at the length level: rewards on incomplete responses can be computed as well, and help steer the trained models towards high-reward responses.

This perspective can be useful to learn to solve problems that can even be a priori combinatorial in nature (by providing some suitable continuous relaxations of such problems using the LLM cross-entropy function to provide a “soft” enforcement of the constraints). The gap between verification and construction is one of the important gaps that can produce new data and games which can later be distilled into the explicit knowledge.

2.3.5. *Constrained optimization in the space of sequences and inverse prompting.* More generally, the log-likelihood (together with suitable contexts) can be used to define various losses on token sequences. A fundamental optimization problem is the prompting problem, which consists of finding the optimal set of tokens that satisfies some fixed constraints (format, length, vocabulary restrictions, ...) and maximizes the log-likelihood when inserted into a fixed context. This family includes some meaningful inverse problems:

Example 5 (Inverse prompting and summarization). Given a fixed text x , we can search for a short sequence of tokens t such that x becomes highly likely when it follows t , i.e.

$$t^* \in \arg \max_{t \in \mathcal{C}} \log \mathbb{P}_{\mathcal{M}}(x | t),$$

where \mathcal{C} encodes the constraints. With appropriate prompting between, before or after t and x , this can be interpreted as a form of summary: t summarizes x because it reduces the surprise of x for the model. This type of implicit knowledge is a direct building block of the subsequent Xent Games.

Example 6 (Common explanations). A variant of the previous example can be obtained when given multiple texts x_1, \dots, x_n . We can search for a sequence of tokens t that makes them jointly more likely, but at the same time relatively unexpected from each of them individually

$$t^* \in \arg \max_{t \in \mathcal{C}} \log \mathbb{P}_{\mathcal{M}}(x_1, \dots, x_n | t) - \alpha \sum_{i=1}^n \log \mathbb{P}_{\mathcal{M}}(t | x_i).$$

With appropriate prompting, t can be understood as an interesting common feature about the texts.

2.3.6. *Anomaly detection.* If a small part of a text has a high cross-entropy loss, it probably deserves attention: this could signal an anomaly, an adversarial or injected segment, or simply a rare but meaningful information in the text that the model cannot well predict. This intuition has been used in practice for e.g. adversarial prompt detection [HWMHS23] and outlier spatial trajectories [MPA24].

2.4. **Explicit vs Implicit Knowledge.** As discussed, the explicit knowledge of an LLM is a subset of its implicit knowledge, since sampling tokens from the model is an algorithm on its probability measure. It is a strict subset, because none of the probabilities in the examples in 2.3 can be computed with a (reasonable) amount of sampling. While some questions like the truth value of a statement can also be answered by generating tokens, having access to the probability allows us to determine the certainty of the model. This access through the model’s probability measure furthermore does not require fine-tuning for instruction-following and is thus well-defined for pre-trained models.

²In the context of self-evaluation, the authors of [KCAHK22] note that “verification improves faster than generation quality in this context”.

Due to its vastness, generality and easy accessibility, we use the implicit knowledge of LLMs as the basis of Xent Games, introduced in the next section. These games will make up the training curriculum of an agent designed for general capabilities.

3. STREAMLINED XENT GAMES

Building on the idea of implicit knowledge outlined in Section 2 above, we briefly outline the ideas behind Xent (Cross-Entropy) Games. Rather than fully specifying the language used to write Xent Games (see [HoEm25] for details), we focus on a core subset of them that we call *Streamlined Xent Games (S-XGs)*, and we provide illustrative examples.

S-XGs form a subspace of the Original Xent Games (O-XGs) introduced in [HoEm25], which focuses on benchmarking. By contrast, S-XGs are designed to make training more streamlined (i.e. more parallelizable, easier to optimize, and with extra differentiable structures [HERG26, HAGER26]) and therefore suitable for training models to play them. At the same time, it is reasonable to expect that the set of skills developed by S-XGs is the same as that developed by the space of O-XGs (see Section 3.3.2 below).

3.1. Definition and Runtime. Informally, Xent Games are text-based games played by some main players \mathcal{M} , with one or more LLMs used as “world models” (i.e. providing the environment dynamics and rewards). The whole space of games is endowed with some meta-data, corresponding to the specification of the models involved, which are each attached to a *model name* (i.e. they link variable names used in the game codes to concrete model checkpoints). In the simplest variant (on which we focus here), one of these is the main player model, which is the one under cognitive training; the other models are frozen and they play the roles of judge, data stream models, or NPCs (opponents / cooperators).

The game state consists of a collection of string registers that are updated according to basic string operations (Section 3.1.1), by using inputs from data streams, and by writing players’ outputs.

The players receive as input strings read from the registers and produce output strings that are written back into the registers. In addition, they receive rewards based on signed cross-entropies of token strings stored in the registers (Section 3.1.2).

3.1.1. Token String Space. As text-based games, Xent Games run on a space of strings, endowed with a small set of basic string operations. In the framework of O-XGs, the allowed operations are: ‘cat’ (i.e. string concatenations at the character level) and ‘cut’ operations (i.e. splits between what comes before a first occurrence of a string and what comes after). In S-XGs, the situation is simpler:

- Strings each have a token length that can change over the course of a game.
- Concatenations and cuts are replaced by copy operations made at the token level (rather than the character level) that copy part of a token string into another. The append operation extends the length of a string up to the maximal length, the cut operation moves tokens from one string to another until the second string’s length is reached.

3.1.2. Xents and Xent Sums. At the heart of Xent Games are cross-entropies of strings $\text{xent}_{\mathcal{J}}$ computed by a (judge) model \mathcal{J} . If $x = x_1, \dots, x_m$ and $y = y_1, \dots, y_m$ are token strings (read from a token string register), we define

$$\text{xent}_{\mathcal{J}}(x|y) = -\log \mathbb{P}_{\mathcal{J}}(x|y) = -\sum_{i=1}^n \log \mathbb{P}_{\mathcal{J}}(x_i|y, x_{<i}).$$

Informally, $\text{xent}_{\mathcal{J}}(x|y)$ measures how “surprised” the autoregressive model \mathcal{J} is to see x after seeing y . We denote also $\text{xent}_{\mathcal{J}}(x) = \text{xent}_{\mathcal{J}}(x|\emptyset)$ where \emptyset is the empty string.

The (judge) model \mathcal{J} can be the main player model \mathcal{M} itself (as in, e.g., pre-training games; see below), or a fixed pre-trained or instruct model. Depending on whether $\mathcal{J} = \mathcal{M}$ or $\mathcal{J} \neq \mathcal{M}$, the goal may be to improve the model used for the cross-entropy computation, or simply to optimize with respect to this cross-entropy defined by a fixed \mathcal{J} .

Remark 7. In practice, for S-XGs, it can be useful to clip the $\mathbb{P}_{\mathcal{J}}(x_i|y, x_{<i})$ from below by a small amount, which clips the per-token loss $-\log \mathbb{P}_{\mathcal{J}}(x_i|y, x_{<i})$ from above. This prevents players from exploiting very low probabilities of models to achieve high rewards: these log-probabilities are typically noisy and not meaningful (see the well-posedness criterion in [HoEm25]).

For a family of token strings $\{(x_j, y_j)\}_j$ and a family of models $(\mathcal{J}_j)_j$ a signed *xent sum* is an expression of the form

$$\sum_j \sigma_j \text{xent}_{\mathcal{J}_j}(x_j|y_j), \quad \sigma_j \in \{\pm 1\}$$

Xent sums can be used in two ways:

- as rewards given to the players;
- as skewed rewards (i.e. soft constraints, see next Section 3.1.3) given to the players.

3.1.3. *S-XG: Definition and Basic Runtime Instruction.* An S-XG consists of a sequence of moves of four types:

- *assign*: modify a string register using the token-string operations described in Section 3.1.1.
- *elicit*: request from a model a string of length n tokens.
- *reward*: reward a player based on a signed xent sum (evaluated on token-string registers).
- *ensure*: provide a smooth constraint on a player by imposing a soft positivity constraint $S \geq 0$ on a signed xent sum S (evaluated on token-string registers). Concretely, for a fixed $\lambda > 1$, the model receives reward S/λ if $S \geq 0$ or λS if $S < 0$ (see Remark 8 below).

Remark 8. The S-XG *ensure* constraints follow the same “adversarial multiplier” logic as in O-XG [HoEm25], but with a fixed finite multiplier $\lambda > 1$ in case of violation (and a small nonzero multiplier $1/\lambda$ in case of fulfillment). Natural examples of *ensure* constraints are true/false statements based on the implicit knowledge (as in Section 2.3.1). Another example of an *ensure* constraint is based on the signed xent sum

$$\text{xent}_{\mathcal{J}}(s2|s1 + \text{"comes after"}) - \text{xent}_{\mathcal{J}}(s1|s2 + \text{"comes after"}),$$

which softly enforces that $s2$ is more likely to follow $s1$ than the reverse.

3.1.4. *Original Xent Game Language (O-XGL).* The simple principles behind O-XGs make them naturally suited to being expressed in a domain-specific language, which we call *O-XGL* (Original Xent Game Language). The design is close in spirit to an assembly language: each line of code corresponds to one instruction of the game logic. Any program made of valid lines of code is valid (there are no conditional jumps), which allows for e.g. a batched execution for training. Each line corresponds to a basic instruction of the types listed in Section 3.1.3 above. In the O-XGL specification [HoEm25], a number of instructions made to simplify the writing and reading of Xent Games by humans are added; however, they do not alter the space of games being considered, and for simplicity, we will omit those here.

The inputs to the *reward* and *ensure* statements are signed xent sums. In O-XGL, these instructions take three ingredients as inputs: the judge, the target string, and the (possibly empty) prefix strings.

3.1.5. *Streamlined Xent Game Language (S-XGL).* As explained above, for the purpose of cognitive training, we rely on streamlined Xent Games (S-XGs). This motivates the S-XGL (Streamlined Xent Game Language) specification. S-XGL is a minimalistic, assembly-like language, which only consists of two basic binary operators, \ll and \gg , acting on two different types of objects:

- Models (main player, judges, data streams, NPCs). Each model has a context register and a score register, used to accumulate xent-based scores and to implement the *ensure* moves.
- Token strings with individual lengths and a global fixed maximal length, implementing the token-space operations described above (see Section 3.1.1).

For example, given a model \mathcal{M} and a string register \mathcal{S} :

- $\mathcal{S} \ll \mathcal{M}$ samples from \mathcal{M} and appends the resulting tokens to \mathcal{S} until the length of \mathcal{S} has doubled or the maximal length is reached,
- $\mathcal{M} \ll \mathcal{S}$ appends \mathcal{S} to the context of \mathcal{M} ,
- $\mathcal{S} \gg \mathcal{M}$ adds $\text{xent}_{\mathcal{M}}(\mathcal{S})$ to the xent accumulator of \mathcal{M} , and $\mathcal{M} \gg \mathcal{S}$ subtracts it from the xent accumulator,
- $\mathcal{M}_\ell \ll \mathcal{M}_r$ rewards \mathcal{M}_ℓ with the value in the xent accumulator of \mathcal{M}_r , and clears that xent accumulator.

Since the two binary operators can act on an ordered pair of objects (of two possible types), this results in $2 \times 2 \times 2 = 8$ types of operations. Together with the empty line operation, which clears all strings, context registers and score registers, this makes $8 + 1$ operations, which constitute all of S-XGL instructions. All other types of lines are valid code (with no execution value), which can be loaded into token strings using S-XGL instructions. A detailed description of S-XGL is given in Appendix A, and an implementation can be found on the <https://www.github.com/xentlabs/s-xgl> repository.

3.2. Examples of Xent Games. In this subsection, we review a number of Xent Games as a means of illustrating the potential of the space. While many Xent Games can be hand-designed for various purposes (see [HoEm25] for a few), our thesis is that the means towards achieving general capabilities is to rely on the automated design of such games.

3.2.1. Pre-Training Game. The classical pre-training objective, i.e. the minimization of $\text{xent}_{\mathcal{M}}$ is a canonical Xent Game. This game is in some sense the simplest: there is no move elicited from \mathcal{M} , just a reward $-\text{xent}_{\mathcal{M}}(x)$ where x is a random string loaded from a data model.

Variants of the basic pre-training objective, as e.g. the multi-token prediction objective [Dee24] (e.g. asking a model to predict the value for 10 tokens ahead, without seeing the 9 ones coming before) can naturally be represented as Xent Games.

3.2.2. Reinforcement Learning as Pre-Training (RLP). In [HAPCC25], a certain game, called RLP, is considered (inspired by a similar game [DDSW25], called RPT). The RLP task is in fact a Xent Game. Changing the notation compared to the paper to fit the Xent Game description, and removing a number of details that are not relevant to our discussion, the game is the following: given a *game map* $x_{<t}, x_t$ taken from a dataset (i.e. where x_t follows the tokens of $x_{<t}$), the game elicits a string c from the player \mathcal{M} to improve the prediction of x_t given $x_{<t}$ and c : the reward of \mathcal{M} is $-\text{xent}_{\mathcal{M}}(x_t|x_{<t}, c)$. This is an interesting example of a game where the judge and the player are the same model.

3.2.3. Distillation and Self-Distillation Games. Distillation and self-distillation can both be viewed as instances of Xent Games.

In online distillation [LTML25, AVZSB24], we have a judge (teacher) model \mathcal{J} that we want to distill into a player (student) model \mathcal{M} . The associated Xent Game uses the contrastive objective across models (Section 2.3.3). Concretely, we provide to \mathcal{M} a context x (possibly void), and we elicit some answer c from \mathcal{M} . The reward for \mathcal{M} is then $\text{xent}_{\mathcal{M}}(c|x) - \text{xent}_{\mathcal{J}}(c|x)$. Maximizing this reward corresponds to minimizing a reverse-KL objective between the (student) model \mathcal{M} and the (teacher) model \mathcal{J} .

In the articles [HLBBG26, SDHA26], the authors do not consider the contrastive objective across models but the counterfactual thinking one (Section 2.3.2), where the extra information comes either from an exact proof, or feedback derived from previous attempts. A simplified Xent Game associated with [HLBBG26] can be described as follows. The main player \mathcal{M} is cloned to obtain a judge \mathcal{J} used to provide reward. We provide x to \mathcal{M} , and elicit a continuation c from \mathcal{M} . Next, conditioned on x, c we prompt \mathcal{J} for feedback and obtain f . The player then receives the reward $\text{xent}_{\mathcal{M}}(c|x) - \text{xent}_{\mathcal{J}}(c|x, f)$.

3.2.4. Prompt Games. Prompt games are our first concrete illustration of the idea of using implicit knowledge to improve the explicit behavior. Inspired by Example 5, the basic reverse prompting game is as follows. Consider a fixed judge \mathcal{J} ; a *game map* consists of a text s . The player must find a prefix t such that, according to \mathcal{J} , s becomes likely after t : we reward the player with $-\text{xent}_{\mathcal{J}}(s|t)$. This game is related to a number of tasks, including creative summarization or jailbreaking. Generalization follows from Example 6 by taking as game map a set of texts s_1, \dots, s_n : the player must find t that makes them jointly more likely (for \mathcal{J}), while remaining relatively unexpected from each of them individually. Depending on the implementation details (constraints, prompting, regularization...), variants of this inverse prompting game have different interpretations, as discussed below.

- Creative summarization. If we add constraints that prevent token copying (length constraint, no common words, ...), a regularization term such as $-\text{xent}_{\mathcal{J}}(t)/2$ (which is equivalent to taking $-2\text{xent}_{\mathcal{J}}(s|t) - \text{xent}_{\mathcal{J}}(t)$) to obtain a well-formed text t , and insert some sandwich prompt between s and t , then the optimal prefix t can be considered as a summary of s .
- Prompt injection. Prompt injection tasks can be formulated as Xent Games. A simple case is the static objective: find a prompt (or injected data) that causes a model \mathcal{J} to generate a fixed output s , independent of the user’s instructions i and data d . To evaluate the gap between the response generated by the LLM \mathcal{J} and the target answer s , the authors of [LYZZX24] propose using a xent loss, e.g. $\text{xent}_{\mathcal{J}}(s|i, d, x)$. In the corresponding Xent Game, another model (the attacker) proposes x (conditioning on (i, d) or not), to maximize the expected reward over a training set of instructions-data pairs.
- Defensive variants. Defensive methods such as DataSentinel [LJSG25] can be viewed as a two-player Xent Game, whose goal is to prevent injected tasks. Two models, an attacker \mathcal{A} and a detector \mathcal{D} ,

compete; this leads to a game-theoretic min-max optimization problem. The attacker outputs contaminated prompts or data (given some instruction and data to process) so that a judge LLM \mathcal{J} executes an injected task instead of the intended task, while also avoiding to be detected. The detector receives a fixed detection instruction s_d and the (possibly contaminated) prompt/data and outputs a secret detection string which serves as marker. This marker should be present when the input is clean, and absent from the detector’s output if contaminated. In this setup, all losses are cross-entropy terms. For example, detection is quantified by the cross entropy $\text{xent}_{\mathcal{G}}(k \mid s_d, x)$ of the marker k given the detection instruction s_d concatenated with the potentially contaminated data x .

3.2.5. Approximate Verifiable Reward Games. When verification and scoring is easier than construction, a judge model can be used to enforce legality of player moves and score plausibility. As explained in Section 2.3.4, such a well-trained judge, even if imperfect, can provide a continuous signal: reward in a verifiable task is thus turned into continuous feedback that can drive parameter optimization.

One may worry that the player will discover ways to exploit the judge biases rather than solve the intended verifiable task. While it is always possible to add an external reward for truly solving the problem, the Xent Games usage relies on transfer across games: if reward hacking is too easy, we conjecture that the model will not learn new generalizable capabilities, and such a game will be filtered out when games are selected in the curriculum.

Examples of such verifiable reward games include chess and mathematical proofs (as discussed in [HoEm25]):

- **Chess.** Games like chess can be used to define Xent Games, based on a judge model which is strong enough to 1. recognize legal moves and 2. evaluate positions. A game map in this case is a chess position. Two players alternately output moves, ensure functions enforce each move’s validity. After a fixed number of steps, or at termination, the score is defined from the judge’s cross-entropy difference between statements such as “white is winning” and “black is winning”.
- **Mathematical proofs.** Let us again assume that a judge model \mathcal{J} is strong enough to 1. estimate the correctness of short formal proofs (for example, in Lean) and 2. evaluate the plausibility of mathematical proof sketches written in English. We can adapt the debate protocol of [CGHCL21, CGHCL24], where two agents explore a tree of proof sketches: the prover wants to maximize the plausibility of the proof by adding details, while the skeptic wants to minimize it by asking for details. The economic rewards of the SPRIG protocol are then replaced by xent terms using the judge model, defining a zero-sum Xent Game.

3.3. Xent Game Space Properties. As argued in Section 3.2, the space of Xent Games \mathcal{G} forms a vast family, eliciting numerous skills highlighted in other papers. In this subsection, we argue that this space possesses a number of good properties that make it suitable for automated exploration towards building generally capable agents.

3.3.1. Closure under Axioms. The first natural property of the Xent Game space is its closure under natural axioms. In [HoEm25], the O-XG space can be shown to naturally emerge from a small collection of game-theoretic and compositional axioms. The S-XG space can be characterized very similarly:

Claim 9. If a space of text games on token string space (with the operations described in Section 3.1.1) contains the reverse prompt game (Section 3.2.4) and is stable under compositionality (we can append the instructions of one game to another), zero-summing (a player’s loss can be a player’s gain and vice versa) and adversarial rescalings (we can turn *reward* statements into *ensure* ones), then it must contain the space of S-XG.

Remark 10. In [HoEm25], adversarial rescaling is formulated in terms of unbounded multipliers, leading to hard constraints; for the sake of training models (rather than evaluating them), it is better to work with soft constraints (with bounded multipliers); the principle is exactly the same. Similarly, the set of allowed operations on the string space in S-XG is slightly different from that of O-XG, but the principle is exactly the same, and if we use the operations described in Section 3.1.1, we obtain the Xent Games described in Section 3.1.3.

3.3.2. Web-of-Games Assumption and Curriculum Universality. A second important conjectural property of the Xent Game space outlined in [HoEm25] concerns the transfer value: informally, the idea is that it is relatively easy, given a set of games, to “discover” new, related games, allowing one to navigate in the space of games, and to grow the skills of a model. In the context of curriculum learning, this motivates the following approximation claim about curricula made of general verifiable reward environments (i.e that are not necessarily Xent Games):

Claim 11. Let E be a skill evaluation. Suppose there exists a curriculum of verifiable reward environments/tasks $\mathbb{G}_1, \dots, \mathbb{G}_n$ such that training on this curriculum brings any model \mathcal{M} in a family \mathfrak{M} to a target skill level $\Lambda_{\mathcal{M}}$ on E . Then, for any $\epsilon > 0$, there exists

- a judge model $\mathcal{M}_{\mathcal{J}}$,
- a curriculum of Xent Games G_1, \dots, G_m of Xent Games that is algorithmically computable from $\mathbb{G}_1, \dots, \mathbb{G}_n$,

such that training any $\mathcal{M} \in \mathfrak{M}$ on G_1, \dots, G_m achieves skill level at least $\Lambda_{\mathcal{M}} - \epsilon$ on E . Moreover, suppose that the curriculum $\mathbb{G}_1, \dots, \mathbb{G}_n$ is computed by a greedy optimization process, then the curriculum of Xent Games G_1, \dots, G_m can also be computed using a greedy optimization process.

Taken at face value, this claim is not very surprising: if (for instance), we have an environment with binary outputs we can demand that $\mathcal{M}_{\mathcal{J}}$ be strong enough to evaluate the games of $\mathbb{G}_1, \dots, \mathbb{G}_n$ and approximate (naively) the scores by looking at the xent of “the output of \mathbb{G}_j is 0” versus the xent of “the output of \mathbb{G}_j is 1”. This is similar to the ideas of 3.2.5.

Similarly to the universality theorems for neural networks, the means to justify Claim 11 are not necessarily directly informative of practical/relevant uses of Xent Games; rather, the statement’s point is that “nothing is missing” from the space.

In Section 4, we leverage these assumptions to formulate cognitive training as a meta-sampling process over the space of Xent Games.

3.4. Game Training and Transfer. The purpose of Xent Games is to provide a suitable environment for robust learning. In this subsection, we introduce the key notion of transfer value between games, which quantifies how training under a scheme Φ on one game affects the performance on another.

3.4.1. Training Scheme. We consider a fixed training scheme Φ , which defines, for any Xent Game G , a map $\mathcal{M} \mapsto \Phi_G \mathcal{M}$ which returns the model obtained by training \mathcal{M} on G for *one run* of the game. A game with multiple training steps should be thought of in our framework as being made of many concatenated copies of a game $G \oplus G \oplus \dots \oplus G$.

Remark 12. For the sake of training, it can be considered that training steps are performed when a clearing “end-of-game” instruction is called.

We assume that the training scheme Φ is invariant by rescaling (if αG denotes the games obtained by multiplying all scores by α , then training is identical and we obtain $\Phi_{\alpha G} = \Phi_G$; in other words, the training is invariant by a rescaling of units of score). This is the case for GRPO-style training based on centered and normalized rewards, as well as for algorithms such as [HERG26, HAGER26].

Remark 13. Importantly: the space of Xent Games is not suitable to perform arbitrary rescalings of game scores by $\alpha > 0$. We *could theoretically allow* the S-XGL language to support a global rescaling of scores for each game, but this would not change anything about what can be achieved in terms of model training (as a result of the scale invariance of the training). However this theoretical possibility (which reflects a training algorithm) is important for our derivation (see Section 4.3)

3.4.2. Transfer Value. A central quantity in our approach is the transfer value between games. Informally, the transfer value $\mathcal{T}_G^{\mathcal{M}}(H)$ encodes “how much in expectation does training on G teaches \mathcal{M} about how to play H ”. We denote the expected score of \mathcal{M} on H by

$$S_{\mathcal{M}}(H) := \mathbb{E}[\text{Score}_H(\mathcal{M})].$$

The transfer value is defined as follows.

Definition 14. For two games G and H , the transfer value $\mathcal{T}_G^{\mathcal{M}}(H)$ from G to H for \mathcal{M} is

$$\mathcal{T}_G^{\mathcal{M}}(H) := S_{\Phi_G[\mathcal{M}]}(H) - S_{\mathcal{M}}(H).$$

More generally, if E is an external evaluation and G is a Xent Game, we define

$$\mathcal{T}_G^{\mathcal{M}}(E) := S_{\Phi_G[\mathcal{M}]}(E) - S_{\mathcal{M}}(E),$$

where $S(E)$ is the expected reward on E .

Remark 15. While these notions are well-posed theoretically, estimating them in practice may require a large number of samples or other techniques (such as the upcoming [HERG26, HAGER26]).

A useful scaling result is the following:

Remark 16. For games G, H we have the scaling relations:

- $\mathcal{T}_G^{\mathcal{M}}(H \oplus H) = 2\mathcal{T}_G^{\mathcal{M}}(H)$: this is as we are purely evaluating on H (twice) and the model is not learning between steps.
- For a small game G , we have $\mathcal{T}_{G \oplus G}^{\mathcal{M}}(H) \approx 2\mathcal{T}_G^{\mathcal{M}}(H)$; training twice on it yields (in expectation) approximately twice the increase of performance on H (this is only an approximation, as if we keep repeating a game enough, we may see diminishing returns).

In Section 4 below, transfer value will be key to:

- estimating the internal relevance of candidate new games;
- estimating the novelty brought by new games compared to previously selected games.

3.4.3. Positive Correlation. A fundamental assumption upon which cognitive learning rests is that we can work with games that allow us to grow a curriculum constructively in a monotone way, i.e. without “needing to make a step back for every two steps forward”.

For two Xent Games $G_1, G_2 \in \mathcal{G}$, we say that they are positively correlated (relative to a model \mathcal{M}) if $\mathcal{T}_{G_i}^{\mathcal{M}}(G_j) > 0$ for $i, j \in \{1, 2\}$: informally speaking, this simply means that learning one doesn’t make \mathcal{M} worse at the other. To grow a curriculum of games in \mathcal{G} , we want to be able to pick at any time (and then optimize with respect to the meta-objective, see Section 4) new games that are positively correlated with respect to all the past games of the curriculum, as measured at appropriate times (see also Section 4.3 below for a justification of the times): if we have built a curriculum G_0, \dots, G_{k-1} yielding the models $\mathcal{M}_1, \dots, \mathcal{M}_k$, we want to be able to pick a next game G_k such that (*a minima*) for each $j < k$, we have:

$$(3.1) \quad \mathcal{T}_{G_j}^{\mathcal{M}_j}(G_k) > 0,$$

$$(3.2) \quad \mathcal{T}_{G_k}^{\mathcal{M}_k}(G_j) > 0.$$

Definition 17. We define $\mathcal{G}_k^+ \subset \mathcal{G}$ to be the set of Xent Games that are positively correlated to $G_{<k}$.

Remark 18. An undesirable outcome that we want to avoid is that we end up in a *transfer cul-de-sac*, i.e. that we are not able to further pick any new game (even a copy of an old one) satisfying (3.1 and 3.2). In case an otherwise useful game happens to negatively correlate with one (or a few) old games, a solution can be simply to append to it a number of copies of these old games to prevent regression on them (and hence negative correlation).

3.4.4. Relevant Skill Discovery. An important desirable feature of a curriculum of games $(G_k)_{k \geq 0}$, which motivates cognitive training is that of *relevant skill discovery*:

Definition 19. We say that a sequence of games $(G_k)_{k \geq 0}$ achieves *relevant skill discovery* on \mathcal{M} if, when training a sequence of models $(\mathcal{M}_k)_k$, with $\mathcal{M}_0 = \mathcal{M}$ and \mathcal{M}_{k+1} obtained from \mathcal{M}_k by training on G_k , we achieve the following:

- The process keeps discovering new skills, i.e. finding new games that are not completely covered (in the sense of transfer value) by any $G_{<k}$ for a fixed k .
- It keeps growing on already discovered skills, i.e. for any G_j discovered, performance on G_j keeps improving as k increases.

4. COGNITIVE TRAINING

In the previous section, we described the Xent Game Space \mathcal{G} , which we postulate provides a good coverage of the learnable skills of an LLM (see Section 3.3.2). Our goal is now to provide an algorithm to construct a curriculum of games in \mathcal{G} for cognitive training, i.e. to build a “flywheel” that continuously (and autonomously) provides a stream of games that improve the model’s general skills.

Remark 20. In this section, we focus on the training of *a single model* \mathcal{M} , although we may use other frozen (non-trained) models as judges or data streamers. Simultaneously training a family of models can be an interesting way to make cognitive training more efficient, but goes beyond the scope of this paper.

In this section, we describe a meta-objective to gauge the quality of a stream of games, based on elaborations of the evolution-based ideas of [HoEm25]. Given an initial model $\mathcal{M} = \mathcal{M}_0$, the goal is to evolve it into a sequence of models $\mathcal{M}_1, \mathcal{M}_2, \dots$ that are more and more generally capable, where each update $\mathcal{M}_k \mapsto \mathcal{M}_{k+1}$ is obtained by training on a game G_k .

4.1. Cognitive Training as a Greedy Algorithm. The central question is *how to choose the next game in the cognitive training curriculum*. We propose to approach this as a greedy optimization problem, focusing on picking the *next game* to grow the curriculum.

Remark 21. We focus on picking a greedy algorithm, as the task of optimizing a true value function on all curricula is likely completely intractable as the curriculum grows. We postulate that this greedy process is sufficient towards discovering new skills and growing towards general capabilities.

At each step, the goal is to find the best (most valuable) game that allows us to evolve the current model towards higher “cognitive” abilities, i.e. that has acquired “the most valuable new skills” between steps k and $k + 1$.

The *meta-game at stage k* consists in tasking a meta-sampler model \mathcal{M}_S with outputting a Xent Game G_k (in S-XGL): the reward of \mathcal{M}_S is evaluated by a meta-objective $G_k \mapsto \mathcal{O}(G_k)$, which depends on $G_{<k}$ and $\mathcal{M}_{\leq k}$.

The question of determining a principled form for \mathcal{O} is a priori a very underspecified question: assuming that the goal of the meta-game is to ascribe a value to games, this raises the question of whether there is a meta-meta-objective, and so on and so forth... it is not clear that such a problem can be resolved in a constructive fashion, i.e. without relying upon an infinite number of assumptions or parameters.³

4.2. Meta-Objective Formula. In Section 4.3, we show that (perhaps surprisingly) an explicit determination of \mathcal{O} from purely qualitative principles (such as consistency principles) can be obtained, substantially constraining the space of “sound” meta-objectives to an explicit natural choice formula⁴.

Given a history $G_{<k}$ and resulting models $\mathcal{M}_{\leq k}$, and restricting the domain to the positively-correlated games \mathcal{G}_k^+ (see Section 3.4.3), we find

$$(4.1) \quad \mathcal{O} = \frac{qd + bp}{l},$$

where we have

- the *quality* term q is given by $q(H) = \left(\sum_{j<k} \mathcal{T}_H(G_j)\right)^{1-\delta}$ for a diversity hyper-parameter $\delta \in [0, 1]$
- the *diversity* term d is given by $d(H) = \left(\mathcal{T}_H(H) / \sum_{j<k} \mathcal{T}_{G_j}(H)\right)^\delta$,
- the *benchmark* term b is given by $b(H) = \mathcal{T}_H(E)$ for an external benchmark metric term E ,
- the *pressure* term $p \geq 0$ is a hyper-parameter modulating the importance of qd versus b .
- the *length* term l is the raw code length of a game written in S-XGL.

Note that despite a linearly growing number of terms in k for q and d (suggesting quadratic scaling in k to run Cognitive Training), the meta-objective computation is only growing very slowly in k (see Section 4.3.6).

The derivation is presented in Section 4.3 below; in particular, the key result is the derivation of the Internal Meta-Objective Theorem derived in Section 4.3.2.

Remark 22. The hyper-parameters δ and p can in principle depend on the training step; if they do, their formulation is however constrained by similar principles to those used in the derivation below.

4.3. Meta-Objective Principles. In this section, we derive an expression for the meta-objective \mathcal{O} aimed at assessing the value of a new game G_k , given a past curriculum $G_{<k}$ of “old games” used to train the model sequences $\mathcal{M}_{<k}$.

4.3.1. General Principles. The following line of reasoning leads us to a principled expression for \mathcal{O}

- The goal of \mathcal{O} is to measure the value that a new game brings towards training the next model.
- The \mathcal{O} value combines *internal* and *external* relevance terms, denoted by \mathcal{I} and \mathcal{E} respectively.
- Both \mathcal{I} and \mathcal{E} terms measure the “useful work performed” if we train the model \mathcal{M}_k on H .
- The \mathcal{I} and \mathcal{E} terms of a game H are to be normalized by its code length l in number of tokens.
- All terms should be evaluated in an *online* fashion, i.e. not allowing the re-training of archives $\mathcal{M}_{<k}$.
- The \mathcal{E} term depends on the difference in an external benchmark metric E if we train \mathcal{M}_k on H .

³This situation echoes a bit a similar situation in field theory, where a Lagrangian may define an action of govern fields in space-time via e.g. a least-action principles, while there is no least-action principle that a Lagrangian would itself satisfy a priori.

⁴A somewhat analogous situation appears in the so-called bootstrap program in conformal field theory.

We thus obtain the following expression for \mathcal{O} of the form,

$$\mathcal{O} = \frac{\mathcal{I} + \mathcal{E}}{l},$$

in what follows, we will provide a principled derivation for \mathcal{I} and \mathcal{E} .

4.3.2. Internal Meta-Objective Derivation. To provide a principled derivation for the internal relevance \mathcal{I} of the meta-objective \mathcal{O} is quite nontrivial: the problem is a priori heavily underspecified. Assume again that we have obtained \mathcal{M} by training on $G_{<k}$ leading to the sequence $\mathcal{M}_{\leq k}$.

- (1) The term \mathcal{I} depends on the past games $G_{<k}$ via the following transfer values:
 - (a) The “new-to-old” transfer values $\mathcal{T}_H^{\mathcal{M}_k}(G_j)$ for $j < k$: this is the core of the *quality* measure.
 - (b) The “old-to-new” transfer values $\mathcal{T}_{G_j}^{\mathcal{M}_j}(H)$ for $j < k$: this is the core of the *diversity* measure.
 - (c) The “new-to-new” transfer value $\mathcal{T}_H^{\mathcal{M}_k}(H)$: this is a term that is useful for normalization.
- (2) The term \mathcal{I} should be factored as

$$\mathcal{I} = qd,$$

where the *quality* q depends on $\left(\mathcal{T}_H^{\mathcal{M}_k}(G_j)\right)_{j < k}$ and the *diversity* d depends on $\left(\mathcal{T}_{G_j}^{\mathcal{M}_j}(H)\right)_{j < k}$ and on $\mathcal{T}_H^{\mathcal{M}_k}(H)$; using a product (rather than a sum) ensures the optimization of *both* quantities.

- (3) The quality and diversity terms are invariant under *history fusion*: if for $0 < j < k$, we replace the steps

$$\mathcal{M}_{j-1} \xrightarrow{G_{j-1}} \mathcal{M}_j \xrightarrow{G_j} \mathcal{M}_{j+1}$$

by the training step (where [none] is an “idle” game, performing nothing)

$$\mathcal{M}_{j-1} \xrightarrow{G_{j-1} \oplus G_j} \mathcal{M}_{j+1} \xrightarrow{[\text{none}]} \mathcal{M}_{j+1},$$

then values of q and d for any G_k are left unchanged.

- (4) From this, we deduce that q must only depend on the sum of transfer values, i.e. must be of the form

$$q(H) = f_q \left(\sum_{j < k} \mathcal{T}_H^{\mathcal{M}_k}(G_j) \right),$$

$$d(H) = f_d \left(\sum_{j < k} \mathcal{T}_{G_j}^{\mathcal{M}_j}(H), \mathcal{T}_H^{\mathcal{M}_k}(H) \right).$$

- (5) Using the “theoretical rescaling invariance idea” (see Remark 13), we get the following: we must have rescaling invariance for $q(H)$, i.e. $q(\alpha H) = q(H)$ for $\alpha > 0$. From this idea, we also get that $\mathcal{I}(\alpha H) = \mathcal{I}(H)$ (since the value of a game can only depend on its effect for training models and training is invariant under rescaling). From there, we find that $d(\alpha H) = d(H)$ as well. Since $\mathcal{T}_{G_j}^{\mathcal{M}_j}(\alpha H) = \alpha H$ and $\mathcal{T}_{\alpha H}^{\mathcal{M}_k}(\alpha H) = \alpha \mathcal{T}_H^{\mathcal{M}_k}(H)$, we find that f_d must be a function of the *ratio* of its two arguments only, i.e. be of the form

$$d(H) = F_d \left(\frac{\mathcal{T}_H^{\mathcal{M}_k}(H)}{\sum_{j < k} \mathcal{T}_{G_j}^{\mathcal{M}_j}(H)} \right)$$

for an increasing bijective function $F_d : \mathbb{R}_+ \rightarrow \mathbb{R}_+$.

- (6) From Remark 16, we have that $\mathcal{T}_{H \oplus H}^{\mathcal{M}_k}(H \oplus H) \approx 4\mathcal{T}_H^{\mathcal{M}_k}(H)$, and requiring $\mathcal{I}(H \oplus H) \approx 2\mathcal{I}(H)$ (the usefulness of training on a game twice is approximately twice the usefulness of the game), we obtain $qd(H \oplus H) \approx 2qd(H)$ and assuming that f_q and F_d must be homogeneous increasing functions, we find that their powers sum up to 1, yielding

$$\left(\sum_{j < k} \mathcal{T}_H^{\mathcal{M}_k}(G_j) \right)^{1-\delta} \left(\frac{\mathcal{T}_H^{\mathcal{M}_k}(H)}{\sum_{j < k} \mathcal{T}_{G_j}^{\mathcal{M}_j}(H)} \right)^\delta,$$

for some *diversity parameter* $\delta \in [0, 1]$.

From the above, we obtain the following:

Theorem (Internal Meta-Objective Theorem). *A consistent internal meta-objective (according to the principles outlined above) \mathcal{I} must be of the form*

$$\mathcal{I} = qd,$$

where

$$q(H) = \left(\sum_{j < k} \mathcal{T}_H^{\mathcal{M}_k}(G_j) \right)^{1-\delta},$$

$$d(H) = \left(\frac{\mathcal{T}_H(H)}{\sum_{j < k} \mathcal{T}_{G_j}^{\mathcal{M}_j}(H)} \right)^\delta,$$

where $\delta \in [0, 1]$ is a diversity hyper-parameter.

4.3.3. *External Meta-Objective Derivation.* We can apply a similar idea to form a natural external meta-objective:

- (1) Similarly to \mathcal{I} , the term \mathcal{E} should satisfy $\mathcal{E}(H \oplus H) \approx 2\mathcal{E}(H)$.
- (2) As a result, it is natural to postulate that

$$\mathcal{E} = bp,$$

where $b = \mathcal{T}_H^{\mathcal{M}_k}(E)$, with E being the given external benchmark metric, where p is a hyper-parameter.

4.3.4. *On the Hyper-Parameters δ and p .* At every step, the two hyper-parameters $\delta > 0$ and $p \geq 0$ are in principle allowed to vary; a principled derivation of formulae for those (i.e. how they should vary in time) is a priori non-trivial and lies beyond the scope of the present work.

A simple choice for both is to take them as constants or to rely on some scaling with respect to the time $T = \sum_{j < k} l(G_j)$ (note that this expression for T is constrained by history fusion).

4.3.5. *Remarks on the Derivation.* In Sections 4.3.1–4.3.3, a principled derivation for \mathcal{O} is provided. A number of remarks are in order.

- The naming for the *quality* and *diversity* terms q and d comes from the Quality-Diversity line of research in Artificial Life (see e.g. [PSS16, ECMO23]), but it is interesting to see that the derivation is made out of first principles only (as opposed to a modeling attempt).
- The derivation can be viewed as the result of the answer to the question: what are the most natural objects we are able to use, and what is a consistent way to combine them? It is notable that a few natural constraints can narrow down the space of consistent meta-objectives to a few hyper-parameters.
- The derivation is not restricted to Xent Games: this is discussed in Section 4.7.

4.3.6. *Scalability of Meta-Objective Computation.* While we assume access to the whole set of archives \mathcal{M}_j for $j \leq k$, it is notable that to compute the diversity denominator sum $\sum_{j < k} \mathcal{T}_{G_j}^{\mathcal{M}_j}(H)$ only requires us to measure the difference of performance of H on \mathcal{M}_k and \mathcal{M} : summing the Definition 14, we get a telescoping sum, where the internal terms cancel.

Thus to perform Cognitive Training (Section 4.4), we only need to store the *latest archive* of \mathcal{M}_k . Similarly, to estimate the other sum $\sum_{j < k} \mathcal{T}_H^{\mathcal{M}_k}(G_j)$ only obviously requires access to \mathcal{M}_k , but also can be computed exactly in a time *linear in the number of different games used* (and furthermore is easily parallelizable).

4.4. **Cognitive Training Algorithm.** In Section 4.3, we have derived the explicit form for the meta-objective \mathcal{O} that, given an initial segment of the curriculum $G_{<k}$, ascribes value to a new game $H \in \mathcal{G}_k^+$ (the space of Xent Games positively correlated to $G_{<k}$, see Section 3.4.3). From this we first derive a cognitive training step:

Definition 23 (Cognitive Training Step). Given a step k and a maximal length L , and a sequence of games $G_{<k}$ used to train a sequence of models $(\mathcal{M}_k)_{k \geq 0}$ with $\mathcal{M}_0 = \mathcal{M}$ and \mathcal{M}_{k+1} being obtained by training \mathcal{M}_k , optimize the meta-objective value $\mathcal{O}(H)$ on games of length $l \leq L$.

Remark 24. Playing this game is a priori hard, as it involves an optimization on the space of all games. As a result, playing can only be imperfect, and can be hard to learn (see Section 4.6.1).

From there, the cognitive training is defined as follows:

Definition 25. Given a meta-sampler \mathcal{M}_S (possibly with some prompt), the Cognitive Training consists in the iteration of Cognitive Training Steps (see Definition 4.4) played by \mathcal{M}_S .

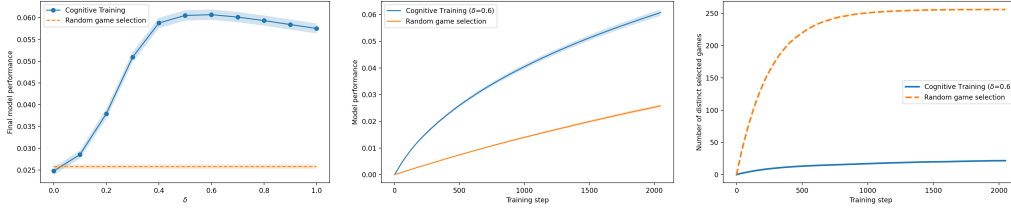


FIGURE 4.1. Left: Geometric mean of game scores as a function of the diversity hyperparameter δ . Middle: geometric mean of game scores during training. Right: number of distinct games selected during training. All values are averaged over 64 independently sampled toy worlds; shaded bands show standard error across worlds.

Remark 26. Note that a priori, the cognitive training algorithm is not a meta-game itself (it has no single objective), being rather the result of iterated play; it is rather a greedy optimization algorithm.

A toy model implementation of Cognitive Training, highlighting the role of the diversity parameter δ , can be found in <https://www.github.com/xentlabs/cognitive-training>.

4.5. Toy Model for Cognitive Training. To give an intuition for the Cognitive Training process, we simulate the greedy optimization of the meta-objective \mathcal{O} on a toy example:

- We assume the existence of a finite list of independent latent skills.
- Each game G is associated with a skill vector that describes how much each skill is required for playing G (and, accordingly, how much playing G improves each skill).

Modeling details are given in Appendix B. Figure 4.1 shows the results of running the algorithm for 2048 steps, comparing it with a baseline that randomly selects a game to train on at each step.

- The left panel shows the dependency of the training on δ : pure quality maximization ($\delta = 0$) performs slightly worse than random choice, whereas any sufficiently large δ leads to a substantial improvement over the baseline.
- The middle panel additionally indicates faster learning with the greedy algorithm throughout the whole training process.
- The right panel shows that the greedy algorithm uses only a small subset of all 256 available games.

4.6. Outlook.

4.6.1. Playing the Meta-Game. The question of playing the meta-game well, i.e. of how to train the meta-sampler \mathcal{M}_S to achieve strong performance on the meta-objective \mathcal{O} is delicate and beyond the scope of this paper: the space of possible Xent Games is enormous and enumerating all games is simply out of the question.

While it is not clear a priori how strong a generalist pre-trained model needs to be to learn to play \mathcal{O} , there is, to the best of our knowledge, no text game with verifiable rewards that is playable by a human that an LLM of current sizes (as of Spring 2026) cannot learn to play at least at a roughly comparable level. A natural setting to allow for the training is to increase the number of samples. This can be achieved by e.g. relying on faster transfer learning estimates to increase sample throughput (such as the upcoming [HERG26, HAGER26]). Another promising way is to train small models \mathcal{M} and use this as a baseline to transfer to larger models.

4.6.2. Open-Ended Skill Discovery. The most significant claim that can be made in the context of cognitive training is that it can lead to an answer to Problem 2, i.e. to *relevant skill discovery* (in the sense of Definition 19 above). The claim is based on the following:

- On the absence of examples of verifiable-reward games that humans can play and that reasonably large language models cannot learn to play.
- On the implausibility that we can achieve very high values of qd on Xent Games having $\min(q, d)$ very small for a nontrivial range of δ -values $\Delta \subset [0, 1]$.

Claim 27. There is a reasonably large model \mathcal{M} and meta-sampler model \mathcal{M}_S such that we have the following: for a nontrivial range $\Delta \subset [0, 1]$ of δ -values, cognitive training with external pressure $p = 0$ yields unbounded relevant skill discovery.

Adding an external benchmark E for which a certain score level λ_E is achievable using a curriculum \mathbb{G} (not necessarily consisting of Xent Games), an analogous question is raised.

Claim 28. Assume $(E, \lambda_E, \mathbb{G})$ is given. There are reasonably large model \mathcal{M} and meta-sampler model \mathcal{M}_S (with access to \mathbb{G} in its prompt), a nontrivial range $\Delta \subset [0, 1]$ of δ -values, and a pressure schedule p , such that cognitive training yields unbounded relevant skill discovery, while converging to λ_E performance on E .

Remark 29. The access to \mathbb{G} for \mathcal{M}_S is not an ideal constraint (which can be voided in natural cases), but it is needed to avoid some situations where learning E involves access to e.g. a cryptographic secret.

4.7. Taking a Step Back. The derivation of the meta-objective obtained in Section 4.3 to get a meta-objective \mathcal{O} for a greedy curriculum is in fact not very specific to Xent Games. The needed ingredient is a concrete framework to define training on games written in a specific language satisfying the following:

- Composition, i.e. any two games G, H can be concatenated to yield a game $G \oplus H$ of length $l(G \oplus H) = l(G) + l(H)$.
- A transfer value definition $\mathcal{T}_G^M(H)$ satisfying:
 - $\mathcal{T}_{G \oplus G}(H) \approx 2\mathcal{T}_G(H)$ and $\mathcal{T}_G(H \oplus H) = 2\mathcal{T}_G(H)$.
 - $\mathcal{T}_{\alpha G}(H) = \mathcal{T}_G(H)$ for $\alpha > 0$ (thus making the explicit support of rescaling unnecessary, see Remark 13).

From this framework, assuming we have a greedy online curriculum of games that bring relevant new skills to the model, we are naturally led to a meta-objective of the form

$$\mathcal{O} = \frac{qd + bp}{l},$$

where q, d, b, p are as above (see Section 4.2).

The derivation’s spirit is the following: we *assume a consistent meta-objective \mathcal{O} exists*, and we find consistency relations it must satisfy, based on the idea that *the only way in which a sequence of games matters is through its effect on model training* (and the only way to judge of this is by judging on its performance on existing games or external metrics). By allowing the space of games to have enough flexibility, we find that this (somehow surprisingly) substantially constrains the shape that the meta-objective must take: if we require that it must not depend on various descriptions of (effectively) the same curriculum, then we find the above explicit formula (leaving only two hyper-parameter degrees of freedom). This line of argument is closely related to the concept of gauge symmetry in theoretical physics.

5. PERSPECTIVES AND CONCLUSION

The main claim of this article is informally the following:

Claim 30. If we believe a greedy curriculum of games to build general capabilities in an LLM can be defined, then it can be replicated by a greedy curriculum of Xent Games that brings relevant new skills (in the sense of Definition 19); and if such a curriculum of Xent Games can be produced, then (under consistency assumptions), its meta-objective must have the shape of the meta-objective \mathcal{O} given by Expression 4.1 in Section 4.2.

Cognitive training based on this meta-objective hence seems to offer a compelling way to build models endowed with general capabilities. It is worth noting that while the space of Xent Games is a compelling framework to implement cognitive training, the principles behind it and the meta-objective derivation can be implemented on any space of games that is sufficiently rich in terms of compositional structure, as discussed in Section 4.7 above.

Acknowledgements. Many insights presented in this paper have emerged from collaborations with Diego Dorn, Vassilis Papadopoulos, Marco Tuccio, Jérémie Wenger, and Nicolas Zlatoff on Large Language Models, with whom key ideas were discussed and investigated.

In addition, the authors would like to thank Emmanuel Abbé, Apoorv Agarwal, Mathieu Alain, Mohammad Asani, Alberto Bietti, Gloria Capano, Rahul Chalamala, Tarun Chitra, Jordan Cotler, Davide Crapis, Marco De Rossi, Mario Geiger, Evgenii Golikov, Alex Graves, Nicola Greco, Bara Hudcová, Arthur Jacot, Niels Linnemann, Tomáš Mikolov, Clément Moulin-Frier, Max Nye, Mihir Patel, João Penedones, David Pfau, Maciej Rudzinski, Stanislav Smirnov, Yi Sun, George Walker, Miles Wang, Jason Wang, Shouqiao Wang, and Matthieu Wyart for interesting discussions, as well as the participants in the Quine seminar and Demeco workshop for insightful questions and remarks.

APPENDIX A: S-XGL SPECIFICATION

In this appendix, we provide a specification of S-XGL (see Section 3.1.5), as implemented on the following repository: <https://www.github.com/xentlabs/s-xgl>. The language specification is deliberately minimal (even more so than that of O-XGL defined in [HoEm25]), consisting of 8+1 instructions. A few preliminary remarks are in order:

- S-XGL is suitable to reward (and thus train) several models; in the cognitive training described above, a single model is trained, and the rewards to other models should just be discarded by the interpreter.
- Some design elements can appear cumbersome to a human reader; the goal is not to make the code particularly human-readable, but to follow the design elements outlined in Section 3.1.5.
- There is no metadata associated with an S-XGL program; S-XGL programs can be seamlessly concatenated.
- Some syntactic sugar can be added to shorten the code without altering the instruction set; we avoid these questions here.

Global Metadata. An S-XGL game space \mathcal{G} consists of the following constant and global (i.e. common to all games $G \in \mathcal{G}$) metadata:

- a fixed token vocabulary \mathcal{V} ;
- a list of models m_u for $u < U$, each based on \mathcal{V} (with $m = m_0$ being the player model under cognitive training);
- a list of token strings s_k for $k < K$ of common maximal length $L \geq 1$.

Code Structure.

- It is understood that any S-XGL program ends with an empty line, which is the clearing instruction, followed by a newline symbol (to enable seamless concatenation of programs); beyond this constraint, any token sequence is valid S-XGL code.
- The interpreter will go line by line through the game code and interpret any line that matches the syntax of an *instruction line* (i.e. that satisfies the syntax of one of the 9 instructions below) as described.
- Otherwise, the code of the line is not interpreted (though it can be used to fill token strings by further lines, as described below; instruction lines can also be used as data to fill token strings).

Set of Variables. Given the global metadata, we define the following set of variables that are allowed to evolve during a game run:

- The token strings s_k for $k < K$: each token string s consists $\text{len}(s)$ tokens in \mathcal{V} , with $0 \leq \text{len}(s) \leq l_{\max}$; at initialization and after clearing, we have $\text{len}(s) = 0$.
- Each model m_u for $u < U$ has a context register and a xent accumulator.

The 8+1 possible expressions.

- $m \ll s$: append the string s to the context register of m .
- $s \ll m$: elicit $\text{len}(s)$ tokens from m (given its current context) and append them to s , stopping upon s reaching the maximal length l_{\max} ; if s is empty, elicit 1 token from m .
- $s \gg m$: compute the xent of s under m (given no context), and add it to the xent accumulator.
- $m \gg s$: compute the xent of s under m (given no context), and subtract it from the xent accumulator.
- $m_\ell \ll m_r$: reward m_ℓ with the score from the xent accumulator of m_r (the judge) and then clear that xent accumulator.
- $m_\ell \gg m_r$: reward m_r with an ensure nonlinearity (see Section 3.1.3) applied to the xent accumulator (the judge) and then clear that xent accumulator.
- $s_\ell \ll s_r$:
 - if $\text{len}(s_r) > 0$: append the tokens of s_r to those of s_ℓ and stop upon s_ℓ reaching maximal length L . In other words, append the first $\min(\text{len}(s_r), l_{\max} - \text{len}(s_\ell))$ tokens from s_r to s_ℓ .
 - if $\text{len}(s_r) = 0$: append the tokenization from the previous code line c to s_ℓ and stop upon s_ℓ reaching maximal length, i.e. append the first $\min(\text{len}(c), l_{\max} - \text{len}(s_\ell))$ tokens of c .
- $s_\ell \gg s_r$: replace s_r with the first $\min(\text{len}(s_\ell), \text{len}(s_r))$ tokens of s_ℓ and remove those tokens from s_ℓ (shifting it to the left and reducing its length by $\min(\text{len}(s_\ell), \text{len}(s_r))$). In particular, if $\ell = r$: clear s_ℓ .
- empty line: clear all strings, context registers and xent accumulators.

APPENDIX B: SIMULATION OF COGNITIVE TRAINING ON A TOY EXAMPLE

We describe the greedy optimization algorithm used to simulate the result of Section 4.5, yielding in particular the results shown in Figure 4.1. To illustrate the dynamics of the algorithm, and in particular its behavior with respect to δ , we use a finite-world toy analogue of the Cognitive Training algorithm.

- The game space \mathcal{G} consists of N_G games, all with the same description length. Each game $G \in \mathcal{G}$ is represented by a non-negative skill vector $(w_{G,s})_{s \in \mathcal{S}} \in \mathbb{R}_+^{N_S}$, where \mathcal{S} is the set of (in practice non-observable) latent skills.
- The game–skill matrix $W = (w_{G,s})$ is sampled as follows: each entry is sampled from a low-mean exponential distribution $\text{Exp}(\mu_{\text{low}})$, and with low probability p_{high} we add a sample from a high-mean exponential distribution $\text{Exp}(\mu_{\text{high}})$. Adding this second component follows the idea that some games strongly rely on certain skills, and that these skills can have a high transfer to some other games.
- A model state is described by a vector of free skills $f = (f_s)_{s \in \mathcal{S}}$ that is mapped to normalized skills through a saturating transformation

$$\text{sk}_S = 1 - (1 + \eta f_s)^{-\alpha_s} \in [0, 1],$$

where $\eta > 0$ is a scale parameter, and $\alpha_s > 0$ is a skill-specific exponent (chosen randomly as μ_{skill} times a squared normal distribution).

- Given a model state f , the score of the model on the game G is defined by

$$S_f(G) = \left(\frac{\sum_{s \in \mathcal{S}} w_{G,s} \text{sk}_s^p}{\sum_{s \in \mathcal{S}} w_{G,s}} \right)^{\frac{1}{p}},$$

which interpolates between the arithmetic mean ($p = 1$) and the weighted geometric mean ($p \rightarrow 0$); in the simulation, we use $p = 1$.

- For a game G , training is modeled by the update $f_s \mapsto f_s + w_{G,s}$, so that the transfer value from G to H is

$$\mathcal{T}_G^f(H) = S_{f+w_G}(H) - S_f(H).$$

- We define the performance of the final model as the geometric mean of the game scores:

$$\left(\prod_{g \in \mathcal{G}} S_f(g) \right)^{\frac{1}{N_G}}.$$

We do not use external benchmarks and a constant description length for the games; the meta-objective thus only consists of the internal relevance term. We apply the Cognitive Training algorithm to construct the curriculum by greedily selecting, at each step, the game H that maximizes the meta-objective for the chosen value of δ . In the example illustrated in Figure 4.1, we run the experiments with $n_{\text{games}} = 256$, $n_{\text{skills}} = 32$, $\mu_{\text{low}} = 100$, $\mu_{\text{high}} = 1$, $\eta^{-1} = \sqrt{n_{\text{games}} n_{\text{skills}}}$, $\mu_{\text{skill}} = 0.1$ for $n_{\text{steps}} = 2048$ steps over $n_{\text{worlds}} = 64$ independently sampled toy worlds.

REFERENCES

- [AVZSB24] R. Agarwal, N. Vieillard, Y. Zhou, P. Stanczyk, S. Ramos, M. Geist, and O. Bachem, On-policy distillation of language models: Learning from self-generated mistakes. In The Twelfth International Conference on Learning Representations, 2024.
- [AAT2023] A.A. Team, J. Bauer, K. Baumli, S. Baveja, F. Behbahani, A. Bhoopchand, N. Bradley-Schmiege, M. Chang, N. Clay, A. Collister, V. Dasagi, L. Gonzalez, K. Gregor, E. Hughes, S. Kashem, M. Loks-Thompson, H. Openshaw, J. Parker-Holder, S. Pathak, N. Perez-Nieves, N. Rakicevic, T. Rocktäschel, Y. Schroecker, J. Sygnowski, K. Tuyls, S. York, A. Zacherl, and L. Zhang, Human-Timescale Adaptation in an Open-Ended Task Space, <https://arxiv.org/abs/2301.07608>
- [Bax00] J. Baxter, A Model of Inductive Bias Learning, *Journal Of Artificial Intelligence Research*, **12**:149–198, 2000, <https://arxiv.org/abs/1106.0245>.
- [BVLFW25] J. Beck, R. Vuorio, E.Z. Liu, Z. Xiong, L. Zintgraf, C. Finn, S. Whiteson, A Tutorial on Meta-Reinforcement Learning, *Foundations and Trends in Machine Learning* **18**(2-3):224–384, <https://arxiv.org/pdf/2301.08028>.
- [AFKKQH22] C. An, J. Feng, K. Lv, L. Kong, X. Qiu, X. Huang, CoNT: Contrastive Neural Text Generation, *Advances in Neural Information Processing Systems* **35** (NeurIPS 2022), <https://arxiv.org/abs/2205.14690>.
- [BaSt07] B. Banerjee and P. Stone, General Game Learning using Knowledge Transfer, Proceedings of the 20th International Joint Conference on Artificial Intelligence, 2007. <https://www.cs.utexas.edu/~ai-lab/pubs/IJCAI07-bikram.pdf>

- [BCK23] J. Blüm, J. Czech, and K. Kersting, AlphaZe*: AlphaZero-like baselines for imperfect information games are surprisingly strong, *Front. Artif. Intell.*, 12 May 2023, Sec. Machine Learning and Artificial Intelligence, Volume 6, <https://doi.org/10.3389/frai.2023.1014561>, 2023
- [CGHCL21] S. Carré, F. Gabriel, C. Hongler, G. Lacerda, and G. Capano, Smart Proofs via Smart Contracts: Succinct and Informative Mathematical Derivations via Decentralized Markets, <https://arxiv.org/abs/2102.03044>.
- [CGHCL24] S. Carré, F. Gabriel, C. Hongler, G. Lacerda, and G. Capano, Smart Proofs via Recursive Information Gathering: Decentralized Refereeing by Smart Contracts, *Distributed Ledger Technologies: Research and Practice*, **3**(1):1–19 <https://dl.acm.org/doi/10.1145/3595298>, 2024.
- [CWWWX23] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, Xing Xie, A Survey on Evaluation of Large Language Models, <https://arxiv.org/abs/2307.03109>
- [CZJGS24] W.-L. Chiang, L. Zheng, Y. Sheng, A.N. Angelopoulos, T. Li, D. Li, H. Zhang, B. Zhu, M. Jordan, J.E. Gonzalez, I. Stoica, Chatbot Arena: An Open Platform for Evaluating LLMs by Human Preference, <https://arxiv.org/abs/2403.04132>.
- [Cho19] F. Chollet, On the Measure of Intelligence, <https://arxiv.org/abs/1911.01547>.
- [CKKLP25] F. Chollet, M. Knoop, G. Kamradt, B. Landers, H. Pinkard, ARC-AGI-2: A New Challenge for Frontier AI Reasoning Systems, <https://arxiv.org/abs/2505.11831>
- [Clu19] J. Clune, AI-GAs: AI-generating algorithms, an alternate paradigm for producing general artificial intelligence, <https://arxiv.org/abs/1905.10985v2>.
- [CKATT18] M.-A. Côté, A. Kádár, X. Yuan, B. Kybartas, T. Barnes, E. Fine, J. Moore, R.Y. Tao, M. Hausknecht, L.E. Asri, M. Adada, W. Tay, and A. Trischler, TextWorld: A Learning Environment for Text-based Games, *Proceedings of the Computer Games Workshop, International Joint Conference on Artificial Intelligence 2018*, <https://arxiv.org/abs/1806.11532>
- [CoTh06] T.M. Cover, J.A. Thomas, *Elements of Information Theory* (2nd Edition), Wiley, 2006
- [DAW24] B.C. Das, M. H. Amini, and Y. Wu, Security and Privacy Challenges of Large Language Models: A Survey, <https://arxiv.org/abs/2402.00888>.
- [Dee24] DeepSeek-V3 Technical Report, DeepSeek-AI, <https://arxiv.org/abs/2412.19437>.
- [DDSW25] Q. Dong, L. Dong, Y. Tang, T. Ye, Y. Sun, Z. Sui, F. Wei, Reinforcement Pre-Training, <https://arxiv.org/pdf/2506.08007>
- [ECMO23] M. Etcheverry, B. W.-C. Chan, C. Moulin-Frier, P.-Y. Oudeyer, Meta-Diversity Search in Complex Systems, A Recipe for Artificial Open-Endedness? <https://arxiv.org/abs/2312.00455>
- [GHWZ16] N. Ghani, J. Hedges, V. Winschel, P. Zahn, Compositional game theory, *Proceedings of the 33rd Annual ACM/IEEE Symposium on Logic in Computer Science. LICS '18*. New York, NY, US: ACM. pp. 472–481. <https://arxiv.org/abs/1603.04641>
- [GGKPW25] A. Gollakota, P. Gopalan, A. Karan, P. Peale, U. Wieder, When does a predictor know its own loss?, <https://arxiv.org/abs/2502.20375>.
- [GBMMK17] A. Graves, M.G. Bellemare, J. Menick, R. Munos, K. Kavukcuoglu, Automated Curriculum Learning for Neural Networks, *Proceedings of the 34th International Conference on Machine Learning*, Sydney, Australia, PMLR 70, 2017. <https://arxiv.org/abs/1704.03003>.
- [HAPCC25] Ali Hatamizadeh, Syeda Nahida Akter, Shrimai Prabhunoye, Jan Kautz, M. Patwary, M. Shoeybi, B. Catanzaro, Y. Choi, RLP: Reinforcement as a Pretraining Objective, <https://arxiv.org/abs/2510.01265>.
- [HiVC93] G. E. Hinton and D. Van Camp, Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory*, pp. 5–13. ACM, 1993. <https://dl.acm.org/doi/10.1145/168304.168306>
- [HoEm25] C. Hongler and A. Emil, Cross-Entropy Games for Language Models: From Implicit Knowledge to General Capability Measures, <https://arxiv.org/abs/2506.06832>.
- [HERG26] C. Hongler, A. Emil, A. Renard, F. Gabriel, Frost Scores and Cross-Entropy Games: Sample-Efficient Training, in preparation.
- [HAGER26] C. Hongler, A. Renard, F. Gabriel, A. Emil, Higher-Order Frost Scoring Schemes, in preparation.
- [HLBBG26] J. Hübotter, F. Lübeck, L. Behric, A. Baumann, M. Bagatella, D. Marta, I. Hakimi, I. Shenfeld, T. K. Buening, C. Guestrin, A. Krause. Reinforcement Learning via Self-Distillation, <https://arxiv.org/pdf/2601.20802>.
- [HWMHS23] Z. Hu, G. Wu, S. Mitra, R. Zhang, T. Sun, H. Huang, V. Swaminathan, Token-Level Adversarial Prompt Detection Based on Perplexity Measures and Contextual Information, <https://arxiv.org/abs/2311.11509>.
- [HDPHR24] E. Hughes, M.D. Dennis, J. Parker-Holder, F. Behbahani, A. Mavalankar, Y. Shi, T. Schaul, T. Rocktäschel, Open-Endedness is Essential for Artificial Superhuman Intelligence, *Proceedings of the 41st International Conference on Machine Learning*, PMLR **235**:20597–20616, 2024. <https://arxiv.org/abs/2406.04268>
- [Hut05] M. Hutter, *Universal algorithmic intelligence: Sequential Decisions based on Algorithmic Probability*. Springer, Berlin, 2005.
- [JTZA26] X. Ji, R. Tutunov, M. Zimmer, H. B. Ammar. Scalable Power Sampling: Unlocking Efficient, Training-Free Reasoning for LLMs via Distribution Sharpening. <https://www.arxiv.org/pdf/2601.21590>.
- [KCAHK22] S. Kadavath, T. Conerly, A. Askell, T. Henighan, D. Drain, E. Perez, N. Schiefer, Z. Hatfield-Dodds, N. Das-Sarma, E. Tran-Johnson, S. Johnston, S. El-Showk, A. Jones, N. Elhage, T. Hume, A. Chen, Y. Bai, S. Bowman, S. Fort, D. Ganguli, D. Hernandez, J. Jacobson, J. Kernion, S. Kravec, L. Lovitt, K. Ndousse, C. Olsson, S. Ringer, D. Amodei, T. Brown, J. Clark, N. Joseph, B. Mann, S. McCandlish, C. Olah, and J. Kaplan, Language Models (Mostly) Know What They Know, <https://arxiv.org/pdf/2207.05221>, 2022

- [KMHBA20] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, D. Amodei, Scaling Laws for Neural Language Models, <https://arxiv.org/abs/2001.08361>.
- [KGRMI22] T. Kojima, S.S. Gu, M. Reid, Y. Matsuo, Y. Iwasawa, Large Language Models are Zero-Shot Reasoners, <https://arxiv.org/abs/2205.11916>
- [Lal2015] K.N. Laland, T. Uller, M.W. Feldman, K. Sterelny, B.B. Müller, A. Moczek, E. Jablonka, and J. Odling-Smee, The extended evolutionary synthesis: its structure, assumptions and predictions, *Proc. R. Soc. B*, **282**: 20151019, 2015.
- [LeHu07] S. Legg, M. Hutter, Universal Intelligence: A Definition of Machine Intelligence, *Minds & Machines*, 17(4):391-444, <https://arxiv.org/abs/0712.3329>, 2007.
- [LeHu06] S. Legg, M. Hutter, A Formal Measure of Machine Intelligence, IDSIA Technical Report 10-06, [arXiv:cs/0605024v1](https://arxiv.org/abs/cs/0605024v1)
- [LeSt10] J. Lehman and K.O. Stanley, Efficiently evolving programs through the search for novelty, *GECCO '10: Proceedings of the 12th annual conference on Genetic and evolutionary computation*, 2010, <https://doi.org/10.1145/1830483.1830638>.
- [LeSt11a] J. Lehman and K.O. Stanley. Abandoning objectives: Evolution through the search for novelty alone. *Evolutionary Computation*, 19(2):189–223, 2011.
- [LeSt11b] J. Lehman and K.O. Stanley. Evolving a diversity of virtual creatures through novelty search and local competition. In *Proceedings of the 13th annual conference on Genetic and evolutionary computation*, pages 211–218. ACM, 2011.
- [LBL23] P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, B. Newman, B. Yuan, B. Yan, C. Zhang, C. Cosgrove, C. D. Manning, C. Ré, D. Acosta-Navas, D. A. Hudson, E. Zelikman, E. Durmus, F. Ladhak, F. Rong, H. Ren, H. Yao, J. Wang, K. Santhanam, L. Orr, L. Zheng, M. Yuksekgonul, M. Suzgun, N. Kim, N. Guha, N. Chatterji, O. Khattab, P. Henderson, Q. Huang, R. Chi, S. M. Xie, S. Santurkar, S. Ganguli, T. Hashimoto, T. Icard, T. Zhang, V. Chaudhary, W. Wang, X. Li, Y. Mai, Y. Zhang, Y. Koreeda, Holistic Evaluation of Language Models, <https://arxiv.org/abs/2211.09110>
- [LHFL22] X.L. Li, A. Holtzman, D. Fried, P. Liang, J. Eisner, T. Hashimoto, L. Zettlemoyer, M. Lewis, Contrastive Decoding: Open-ended Text Generation as Optimization, *Association for Computer Linguistics*, <https://arxiv.org/abs/2210.15097>
- [LKKYW25] B. Liu, C. Jin, S. Kim, W. Yuan., W. Zhao, I. Kulikov, X. Li, S. Sukhbaatar, J. Lanchantin, J. Weston, SPICE: Self-Play In Corpus Environments Improves Reasoning. <https://arxiv.org/pdf/2510.24684>, 2025.
- [LJSG25] Y. Liu, Y. Jia, J. Jia, D. Song, N.Z. Gong, DataSentinel: A Game-Theoretic Detection of Prompt Injection Attacks, 2025 IEEE Symposium on Security and Privacy (SP), 2190-2208, 2025.
- [LYZZX24] X. Liu, Z. Yu, Y. Zhang, N. Zhang, C. Xiao, Automatic and Universal Prompt Injection Attacks against Large Language Models, <https://arxiv.org/abs/2403.04957>.
- [LWLZD24] R. Liu, J. Wei, F. Liu, C. Si, Y. Zhang, J. Rao, S. Zheng, D. Peng, D. Yang, D. Zhou, A. M. Dai, Best Practices and Lessons Learned on Synthetic Data, *COLM 2024*, <https://arxiv.org/abs/2404.07503>
- [LTML25] K. Lu and Thinking Machines Lab, "On-Policy Distillation", *Thinking Machines Lab: Connectionism*, Oct 2025.
- [MPA24] J. K. Mbuya, D. Pfoser, A. Anastasopoulos, Trajectory Anomaly Detection with Language Models, *SIGSPATIAL '24: Proceedings of the 32nd ACM International Conference on Advances in Geographic Information Systems*, 208–219, <https://arxiv.org/abs/2409.15366>.
- [MCV20] Clara Meister, Ryan Cotterell, and Tim Vieira. If beam search is the answer, what was the question? *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 2173–2185. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.emnlp-main.170. <https://aclanthology.org/2020.emnlp-main.170>
- [MKSLH15] V. Mnih, K. Kavukcuoglu, D. Silver, A.A. Rusu, J. Veness, M.G. Bellemare, A. Graves, M. Riedmiller, A.K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, D. Hassabis, Human-level control through deep reinforcement learning, *Nature* 518, 529–533 (2015). <https://doi.org/10.1038/nature14236>
- [MJB15] T. Mikolov, A. Joulin, M. Baroni, A Roadmap towards Machine Intelligence, <https://arxiv.org/abs/1511.08130v2>.
- [MSBLB17] M. Moravčík, M. Schmid, N. Burch, V. Lisý, D. Morrill, N. Bard, T. Davis, K. Waugh, M. Johanson, M. Bowling, DeepStack: Expert-Level Artificial Intelligence in No-Limit Poker, <https://arxiv.org/abs/1701.01724>.
- [OAI2020] T.B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D.M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, D. Amodei, Language Models are Few-Shot Learners, *Advances in Neural Information Processing Systems* **33**:1877–1901, <https://arxiv.org/abs/2005.14165>.
- [PBRW99] L. Page, S. Brin, R. Motwani, and T. Winograd, The PageRank Citation Ranking: Bringing Order to the Web. *Technical Report. Stanford InfoLab*. <http://ilpubs.stanford.edu:8090/422/>
- [PWH24] V. Papadopoulos, J. Wenger, C. Hongler, Arrows of Time for Large Language Models, *International Conference on Machine Learning 2024*. <https://arxiv.org/abs/2401.17505>, 2024.
- [PBS16] E. Parisotto, J.L. Ba, R. Salakhutdinov, Actor-Mimic: Deep Multitask and Transfer Reinforcement Learning, *International Conference on Learning Representations*, 2016, <https://arxiv.org/abs/1511.06342>.
- [PSS16] J.K. Pugh, L.B. Soros, K. O. Stanley, Quality Diversity: A New Frontier for Evolutionary Computation, *Front. Robot. AI*, 3 - 2016 | <https://doi.org/10.3389/frobt.2016.00040>.

- [PCDSC24] Ji-Lun Peng, Sijia Cheng, Egil Diau, Yung-Yu Shih, Po-Heng Chen, Yen-Ting Lin, Yun-Nung Chen, A Survey of Useful LLM Evaluation, <https://arxiv.org/abs/2406.00936v1>
- [Schmi07] J. Schmidhuber. Gödel machines: Fully self-referential optimal universal self-improvers. *Artificial general intelligence*, pages 199–226. Springer, 2007. <https://sferics.idsia.ch/pub/juergen/gmAGI.pdf>
- [Schmi10] J. Schmidhuber. Formal Theory of Creativity, Fun, and Intrinsic Motivation (1990-2010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230-247, 2010. IEEE, <https://people.idsia.ch/~juergen/ieeecreative.pdf>
- [SYCYL24] Chufan Shi, Haoran Yang, Deng Cai, Zhisong Zhang, Yifan Wang, Yujiu Yang, Wai Lam, A Thorough Examination of Decoding Methods in the Era of LLMs, <https://arxiv.org/abs/2402.06925>
- [SYSKZ24] Sanchit Sinha, Yuguang Yue, Victor Soto, Mayank Kulkarni, Jianhua Lu, Aidong Zhang, MAML-en-LLM: Model Agnostic Meta-Training of LLMs for Improved In-Context Learning, International Conference on Knowledge Discovery and Data Mining, 2024, <https://arxiv.org/abs/2405.11446>
- [Sig23] O. Sigaud, G. Baldassarre, C. Colas, S. Doncieux, R. Duro, P.-Y. Oudeyer, N. Perrin-Gilbert, V.G. Santucci, A Definition of Open-Ended Learning Problems for Goal-Conditioned Agents, <https://arxiv.org/abs/2311.00344>.
- [SHSH18] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. Lillicrap, K. Simonyan, and D. Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play, *Science*, 7 Dec 2018, **362**(6419):1140-1144, DOI: 10.1126/science.aar6404, 2018.
- [SDHA26] I. Shenfeld, M. Damani, J. Hübotter, P. Agrawal. Self-Distillation Enables Continual Learning. <https://arxiv.org/pdf/2601.19897>.
- [SWLL24] Y. Song, G. Wang, S. Li, B.Y. Lin, The Good, The Bad, and The Greedy: Evaluation of LLMs Should Not Ignore Non-Determinism, <https://arxiv.org/abs/2407.10457>.
- [StBy19] Felix Stahlberg and Bill Byrne. On NMT search errors and model errors: Cat got your tongue? Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 3356–3362. Association for Computational Linguistics, 2019. doi: 10.18653/v1/D19-1331. URL <https://aclanthology.org/D19-1331>.
- [SSSJ25] Z. Stojanovski, O. Stanley, J. Sharratt, R. Jones, A. Adefoye, Jean Kaddour, Andreas Köpf, REASONING GYM: Reasoning Environments for Reinforcement Learning with Verifiable Rewards, <https://arxiv.org/abs/2505.24760>
- [TKM23] Q. Tan, A. Kazemi, R. Mihalcea, Text-Based Games as a Challenging Benchmark for Large Language Models, International Conference on Learning Representations, TinyPapers, 2023.
- [Tur50] A. M. Turing, Computing Machinery and Intelligence. *Mind* **49**:433-460, 1950.
- [WLCS19] R. Wang, J. Lehman, J. Clune, K. O. Stanley. POET: open-ended coevolution of environments and their optimized solutions. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pp. 142–151, ACM, 2019.
- [WWZ24] Z. Wang, P. Wang, K. Liu, P. Wang, Y. Fu, C.-T. Lu, C.C. Aggarwal, J. Pei, Y. Zhou, A Comprehensive Survey on Data Augmentation, <https://arxiv.org/abs/2405.09591>
- [WXSLD22] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, D. Zhou, Chain-of-Thought Prompting Elicits Reasoning in Large Language Models, Proceedings of NeurIPS 2022, <https://arxiv.org/abs/2201.11903>.
- [Wil92] R.J. Williams, Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning, *Machine Learning*, **8**(3-4):229 - 256, <https://dl.acm.org/doi/10.1007/bf00992696>.
- [XDCGZ24] Z. Xi, Y. Ding, W. Chen, B. Hong, H. Guo, J. Wang, D. Yang, C. Liao, X. Guo, W. He, S. Gao, L. Chen, R. Zheng, Y. Zou, T. Gui, Q. Zhang, X. Qiu, X. Huang, Z. Wu, Y.-G. Jiang, AgentGym: Evolving Large Language Model-based Agents across Diverse Environments, <https://arxiv.org/abs/2406.04151v1>.
- [XFLZZ25(2025)] X. Xing, Z. Fan, J. Lou, G. Li, J. Zhang, D. Zhang, PretrainZero: Reinforcement Active Pretraining, <https://arxiv.org/pdf/2512.03442>
- [ZHLLC25] J. Zhang, S. Hu, C. Lu, R. Lange, J. Clune, Darwin Gödel Machine: Open-Ended Evolution of Self-Improving Agents, <https://arxiv.org/abs/2505.22954>.