

# Learning to Learn: Curriculum Selection via Soft Q-Learning for Sample-Efficient Language Models

Scott Viteri      Damir Vrabac

## 1 Overview

We propose an algorithm that enables language models to learn what data to learn from. Rather than training on randomly sampled data, the model takes actions that select training examples, receiving reward based on how well it currently predicts held-out data. A Q-function, parameterized as a small network reading the language model’s internal representations, learns to predict the long-run value of each candidate. The policy is derived from the Q-function via a Boltzmann distribution, eliminating the need for a separate policy network. The core hypothesis is that this learned curriculum selection will yield significantly better sample efficiency than random or heuristic curricula.

## 2 Technical Approach

### 2.1 Setup

- Language model (GPT-2) with parameters  $\theta$ , hidden dimension  $d = 768$
- Streaming data source, yielding fresh batches of  $N$  candidate context windows each step
- Fixed held-out evaluation set  $D$ , large; a fresh random subset  $\hat{D} \subset D$  of size  $M$  is sampled each step for reward computation
- Q-network  $Q_\phi$  with parameters  $\phi$ , implemented as a 2-layer MLP ( $d \rightarrow 32 \rightarrow 1$ ) reading the base model’s residual stream (see Section 2.3)
- Target Q-network  $Q_{\bar{\phi}}$  with parameters  $\bar{\phi}$ , a Polyak-averaged copy of  $Q_\phi$  used for stable bootstrap targets
- Boltzmann temperature  $\beta > 0$  controlling exploration
- Learning rate  $\alpha$  for the LM gradient step; learning rate  $\eta$  for the Q-network update
- Discount factor  $\gamma \in (0, 1)$  controlling the effective planning horizon
- Polyak averaging rate  $\tau \in (0, 1)$  for the target network update
- Gradient clipping threshold  $G$  (max global norm; applied to both LM and Q-network gradient steps)

### 2.2 State, Action, and Reward

The state  $s_k$  is the current model parameters  $\theta_k$ , which determine the model’s representations and therefore the features available to the Q-network. In practice,  $s_k$  is never represented explicitly—it is accessed implicitly through the base model’s hidden states when candidates are forwarded through it (Section 2.3).

The action  $a_k \in C_k$  is the choice of which training example to train on at step  $k$ , where  $C_k$  is the batch of  $N$  candidates available at that step. Since data is streamed, each step presents a

fresh candidate batch, and each candidate is seen at most once. The Q-network generalises across batches via the base model’s hidden-state representation.

The reward is the held-out log-probability *after* the training update:

$$r_k = \frac{1}{M} \log p_{\theta_{k+1}}(\hat{D}_k), \quad (1)$$

where  $\hat{D}_k$  is a fresh random subset of  $D$  of size  $M$ , and  $\theta_{k+1}$  results from training on  $a_k$ . This measures the model’s absolute predictive fitness as a consequence of the selected action. Resampling  $\hat{D}_k$  each step prevents the Q-function from overfitting to a fixed evaluation subset.

The Q-function optimises the discounted return  $(1 - \gamma) \sum_k \gamma^k r_k$ . The weights  $(1 - \gamma)\gamma^k$  sum to one, so this is the *exponentially weighted average* of the reward sequence: the Q-value at any state equals the typical per-step reward expected over the next  $\sim 1/(1 - \gamma)$  steps under the current policy. Maximising this objective favours curricula that achieve high predictive fitness quickly—each step spent with poor predictions is penalised, with a horizon controlled by  $\gamma$ .

## 2.3 Q-Network Architecture

The Q-network reads the base model’s internal representations to assess candidate value. For each candidate  $x \in C_k$ :

1. Forward  $x$  through the base model with parameters  $\theta_k$ , extracting the last-layer hidden state at the final token position:  $h_x = \text{last\_hidden}(\theta_k, x) \in \mathbb{R}^d$ .
2. Detach  $h_x$  from the computation graph (no gradients flow into  $\theta$ ).
3. Pass through the Q-network:

$$Q_\phi(s_k, x) = W_2 \sigma(W_1 h_x + b_1) + b_2, \quad (2)$$

where  $\phi = \{W_1 \in \mathbb{R}^{32 \times d}, b_1 \in \mathbb{R}^{32}, W_2 \in \mathbb{R}^{1 \times 32}, b_2 \in \mathbb{R}\}$  and  $\sigma$  is ReLU.

The hidden state  $h_x$  implicitly encodes both the model’s current state (through  $\theta_k$ , which determines the representations) and the identity of the candidate (through the input  $x$ ). By detaching  $h_x$ , we ensure that  $\theta$  is updated only by the language modeling objective, while  $\phi$  is updated only by the TD loss. This prevents the co-adaptation instability that arises when value gradients distort the base model’s representations.

## 2.4 Normalised Discounted Soft Bellman Equation

The Q-function satisfies a normalised discounted soft Bellman equation:

$$Q(s_k, a_k) = (1 - \gamma) r_k + \gamma V_{\bar{\phi}}(s_{k+1}), \quad (3)$$

where  $\gamma \in (0, 1)$  is the discount factor and  $V_{\bar{\phi}}$  is the normalised soft value computed from the target network (see below). The  $(1 - \gamma)$  factor on the reward ensures that a constant reward  $r$  contributes  $r$  (not  $r/(1 - \gamma)$ ) to the Q-value fixed point. Combined with the normalised soft value (which subtracts  $\beta \log N$ ), a constant reward gives  $Q^* = r$  exactly. This keeps Q-values on the same scale as the reward and the MLP output in a numerically comfortable range.

The soft value function is normalised against the uniform prior:

$$V(s) = \beta \log \frac{1}{N} \sum_{a'} \exp(Q(s, a')/\beta) = \beta \left[ \log \sum_{a'} \exp(Q(s, a')/\beta) - \log N \right]. \quad (4)$$

The log  $N$  subtraction removes the entropy bonus from the uniform distribution: when all Q-values are equal at  $q$ , we have  $V(s) = q$  rather than  $V(s) = q + \beta \log N$ . This is the *normalised* logsumexp, or equivalently the log-mean-exp. Without it, the entropy bonus  $\beta \log N$  would be amplified by  $\gamma/(1 - \gamma)$  at the Bellman fixed point, inflating Q-values far above the reward scale.

With both normalisations— $(1 - \gamma)$  on the reward and  $-\log N$  in the soft value—a constant reward  $r$  gives  $Q^* = r$  at the fixed point. Q-values live on the same scale as the reward (around  $-3$  to  $-4$  for GPT-2 perplexity), and the Q-network’s output is directly interpretable as the expected per-step reward over the planning horizon. At convergence,  $Q(s, a)$  represents a normalised discounted sum of future rewards from choosing  $a$  in state  $s$  and following the Boltzmann policy thereafter. With  $\gamma = 0.9$ , the effective horizon is  $\sim 10$  steps.

The policy  $\pi(a | s) \propto \exp(Q(s, a)/\beta)$  is unaffected by the  $\log N$  subtraction, since softmax is shift-invariant.

**Target network.** To stabilise the bootstrap target  $V(s_{k+1})$ , we maintain a target network  $Q_{\bar{\phi}}$  that is a Polyak-averaged copy of the online network  $Q_\phi$ :

$$\bar{\phi} \leftarrow \tau \phi + (1 - \tau) \bar{\phi}, \quad \tau \in (0, 1). \quad (5)$$

The soft value in the TD target is computed from  $Q_{\bar{\phi}}$ , not  $Q_\phi$ . This breaks the feedback loop where the network’s own output serves as its training target, preventing the oscillations and divergence that can occur with function approximation. The online network  $Q_\phi$  is still used for action selection (the Boltzmann policy) and for the prediction  $\hat{Q}$  in the TD loss. With  $\tau = 0.01$ , the target network tracks the online network with a lag of  $\sim 1/\tau = 100$  steps.

## 2.5 Algorithm

The key implementation insight is that each step’s forward pass of candidates through the base model serves two purposes: action selection for the *current* step and the bootstrap target for the *previous* step’s Q update. This avoids a redundant second forward pass.

The LM trains on a *policy-weighted mixture* of all  $N$  candidates at every step, rather than on a single sampled action. Define the Boltzmann weights:

$$\pi_i = \frac{\exp(Q_\phi(s_k, x_i)/\beta)}{\sum_{j=1}^N \exp(Q_\phi(s_k, x_j)/\beta)}, \quad x_i \in C_k. \quad (6)$$

The LM loss is the policy-weighted mixture:

$$\mathcal{L}_{\text{mix}}(\theta; C_k) = \sum_{i=1}^N \pi_i \ell(\theta, x_i), \quad (7)$$

where  $\ell(\theta, x_i)$  is the per-token cross-entropy loss on candidate  $x_i$ . This is precisely the *expected gradient* under the Boltzmann policy:

$$\nabla_\theta \mathcal{L}_{\text{mix}} = \mathbb{E}_{x \sim \pi(\cdot | s_k)} [\nabla_\theta \ell(\theta_k, x)], \quad (8)$$

computed exactly rather than estimated via a single Monte Carlo sample. This eliminates the variance from action sampling.

**Q-learning and the mixture LM update.** The Q-network update is unchanged by the mixture formulation: a discrete action  $a_k$  is still sampled from the Boltzmann policy, its hidden state  $h_{\text{prev}}$  is stored, and the standard TD update is applied at the next step. The Q-function continues to satisfy the discounted soft Bellman equation (3). The only difference from a single-action formulation is how the environment *transitions*: the LM parameters are updated via the full mixture rather than a single example.

**Off-policy interpretation.** From the Q-learner’s perspective, the situation is analogous to off-policy learning. The Q-learner observes: “I selected action  $a_k$ , the environment transitioned to state  $s_{k+1}$ , and I received reward  $r_k$ .” It does not need to know that the transition was caused by a mixture-weighted gradient step rather than a single-example step. Q-learning is inherently off-policy—the Bellman backup  $Q(s, a) = (1 - \gamma)r + \gamma V(s')$  is valid regardless of the behavior policy that generated the transition, as long as  $r$  and  $s'$  are the actual observed reward and next state.

Over many steps, the Q-function converges to predict the value of upweighting a given candidate within the mixture—which is exactly the marginal contribution signal needed to set the mixture weights well.

**Effective batch size.** The mixture-weighted update is equivalent to training on a soft batch of  $\sim 1/\|\pi\|_2^2$  effective examples per step (the inverse of the collision probability of the policy distribution). Early in training, when Q-values are similar, this approaches  $N$ ; as the policy sharpens, the effective batch size decreases toward 1.

**Computational cost per step.** One forward pass of  $N$  candidates through the base model (batched, no gradient) for Q-scoring;  $N$  forward-backward passes through the full model for the mixture LM loss (processed in mini-batches with gradient accumulation); one forward pass of  $M$  held-out examples for reward; one MLP forward-backward for the Q update; one MLP forward pass through the target network (negligible). The candidate extraction forward passes dominate; the mixture LM step adds roughly a  $2\times$  factor to total step time compared to a single-example LM step.

**Staleness.** The Q update at step  $k$  uses the target network  $Q_{\bar{\phi}}$  for the bootstrap value  $V_k$ , and the online network  $Q_\phi$  for action selection. The target network lags the online network by  $\sim 1/\tau$  steps, which is by design: this lag stabilises the TD target. The online Q-values  $q_k$  used for action selection have one-step staleness (computed before the  $\phi$  update), which is standard in online Q-learning and negligible given the small MLP learning rate  $\eta$ .

## 2.6 Uniform Mixture Variant (No Q-Weighting)

As an ablation, the flag `--no_q_weighting` replaces the Boltzmann policy weights with uniform weights:

$$\pi_i = \frac{1}{N}, \quad i = 1, \dots, N. \tag{9}$$

The LM loss becomes the unweighted average over all candidates:

$$\mathcal{L}_{\text{uniform}}(\theta; C_k) = \frac{1}{N} \sum_{i=1}^N \ell(\theta, x_i). \tag{10}$$

This decouples the LM training step from the Q-function entirely: the LM trains on an equal mixture of all candidates regardless of their Q-values. The Q-network still trains normally (receiving rewards, computing TD updates), but its output does not influence which data the LM learns from. This serves as a controlled ablation to isolate the effect of Q-weighted curriculum selection on sample efficiency.

## 2.7 Interpretation as Free-Energy Minimisation

The Boltzmann policy is the solution to a free-energy minimisation problem at each state:

$$\pi^*(\cdot | s) = \arg \min_{\pi} \left[ -\mathbb{E}_{x \sim \pi} [Q_\phi(s, x)] + \beta \text{KL}(\pi \| p_0) \right], \quad (11)$$

where  $p_0$  is the uniform distribution over  $C_k$ . Since  $p_0$  is uniform,  $\text{KL}(\pi \| p_0) = \log N - H(\pi)$ , and the objective reduces (up to constants) to:

$$\pi^*(\cdot | s) = \arg \min_{\pi} \left[ -\mathbb{E}_{x \sim \pi} [Q_\phi(s, x)] - \beta H(\pi) \right]. \quad (12)$$

The first term encourages selecting high-value candidates; the second encourages exploration. The temperature  $\beta$  controls the tradeoff.

**ELBO interpretation.** Define a latent variable model where the “evidence” is the event that the current trajectory is optimal, with likelihood  $p(\text{optimal} | x) \propto \exp(Q(s, x)/\beta)$ . The policy  $\pi$  acts as a variational posterior over actions, and the free energy is the negative ELBO:

$$\log p(\text{optimal} | s) \geq \mathbb{E}_{x \sim \pi} \left[ \frac{Q(s, x)}{\beta} \right] - \text{KL}(\pi \| p_0). \quad (13)$$

The Boltzmann policy makes this bound tight. The normalised soft value  $V(s) = \beta[\log \sum_x \exp(Q(s, x)/\beta) - \log N]$ —which appears as the bootstrap term in the TD target—is the log-evidence minus  $\beta \log N$ , measuring how much value the current candidate batch offers above the uniform baseline.

## 3 Evaluation

### 3.1 Primary Metric: Sample Efficiency

The central question is: how many training examples does the model need to reach a given performance level?

We will measure:

- Perplexity on a fixed evaluation set as a function of training examples seen
- Performance on downstream tasks (e.g., MMLU, HellaSwag) as a function of training examples

### 3.2 Baselines

- **Random curriculum:** Uniform sampling from candidates
- **Loss-based curriculum:** Prioritise high-loss examples
- **Uncertainty-based curriculum:** Prioritise examples with high model uncertainty
- **Competence-based curriculum:** Examples ordered by difficulty (requires difficulty labels)

### 3.3 Analysis

Beyond aggregate metrics, we will examine:

- **Curriculum structure:** Does the learned curriculum exhibit interpretable patterns? Developmental stages? Topic clustering?
- **Exploration dynamics:** How does the policy entropy evolve over training? Does the temperature  $\beta$  produce reasonable exploration?
- **Q-function interpretability:** Which candidates receive high Q-values at different stages of training? Does the Q-function learn to identify examples that are valuable given the model’s current state?
- **Q-value dynamics:** How do Q-values evolve over training? Convergence of Q-values indicates the value function has stabilised; the spread (standard deviation) reflects how strongly the Q-function discriminates between candidates.

## 4 Relation to Prior Work

### 4.1 Curriculum Learning

Graves et al. (2017) use multi-armed bandits to select tasks; Jiang et al. (2018) train a separate “mentor” network to weight examples. Our approach differs by using the language model’s own representations as features for the Q-function and applying soft Q-learning with temporal credit assignment, rather than treating selection as a stateless bandit problem.

### 4.2 Meta-Learning

MAML (Finn et al., 2017) learns initializations for fast adaptation. We share the computational structure—reasoning about the effect of gradient updates—but learn *what* to train on rather than *where* to start.

### 4.3 RL as Inference

The control-as-inference framework (Levine, 2018; Rawlik et al., 2013) casts optimal control as probabilistic inference, with the Q-function serving as an energy function and the optimal policy as the corresponding Boltzmann distribution. Our algorithm instantiates this framework in the curriculum selection setting: the Q-network provides unnormalized log-probabilities over actions, and the Boltzmann policy is derived analytically without requiring a separate policy network. This is the soft Q-learning approach of Haarnoja et al. (2017), applied to curriculum selection with a standard discount factor for Bellman contraction.

### 4.4 Soft Q-Learning

Soft Q-learning (Haarnoja et al., 2017) augments the standard Bellman equation with an entropy bonus, yielding a Boltzmann policy without requiring a separate policy network. SAC (Haarnoja et al., 2018) extends this to continuous actions with a learned temperature. Our setting is simpler: the action space is a finite set of  $N$  candidates, so the Boltzmann policy and soft value can be computed exactly via softmax and logsumexp. Following standard practice in deep Q-learning (Mnih et al., 2015; Lillicrap et al., 2016), we use a Polyak-averaged target network for the bootstrap value, normalise the reward by  $(1 - \gamma)$ , and subtract  $\log N$  from the soft value to keep Q-values on the same scale as the reward.

## 4.5 Intrinsic Motivation

Our objective relates to compression progress (Schmidhuber) and active inference, but avoids the memory requirements of the former and the dark room problem of the latter by using a fixed held-out set as a proxy for predictive success. The absolute log-probability reward (1) provides a direct measure of predictive fitness, connecting to the epistemic homeostasis motivation: the agent is penalised for *being* in a state of poor prediction, not merely for failing to improve.

## 4.6 Connection to Prior Work

This proposal builds directly on Markovian Transformers for Informative Language Modeling (<https://arxiv.org/abs/2404.18988>), which introduces a framework for training language models with RL to produce causally load-bearing Chain-of-Thought reasoning. Both projects share a common structure: use RL to learn intermediate representations that improve prediction on held-out data. In the Markovian Transformers work, the learned object is a CoT:

$$\text{Question} \rightarrow \text{CoT} \rightarrow \text{Answer}$$

In the current proposal, the learned object is a curriculum:

$$\text{Model State} \rightarrow \text{Selected Data} \rightarrow \text{Improved Predictions}$$

The Markovian Transformers work demonstrates that this general approach is tractable and yields large gains on reasoning benchmarks (e.g., GSM8K: 19.6% → 57.1%). The current proposal extends this framework from learning *what to say* to learning *what to study*.

## 5 Broader Motivation

Language models trained to predict text develop remarkable capabilities from a simple objective. Yet they require extensive post-training to behave agentically and arguably lack a kind of global coherence. One hypothesis: the training process is purely observational—the model never takes actions that affect what it observes.

This project is a stepping stone toward studying whether learned curriculum selection produces qualitatively different agents. The immediate goal is demonstrating sample efficiency gains. The longer-term question is whether controlling one’s own learning process contributes to the coherence and agency that current models seem to lack.

### 5.1 Long-Term Direction: Formalizing Homeostasis

Biological agents are shaped by survival pressures. Hunger, pain, and fatigue are not arbitrary reward signals—they are tied to the organism’s continued existence. Current approaches to intrinsic motivation (curiosity, empowerment, compression progress) capture aspects of adaptive behavior but lack this grounding in self-preservation.

A long-term goal of this research program is to formalize homeostasis and survival into a simple, biologically plausible metric that could serve as a foundation for intrinsic motivation in artificial systems. The current project—learning to select data that improves future prediction—is a minimal step in this direction: the agent takes actions that maintain its predictive capacity, a kind of epistemic homeostasis. The absolute reward formulation makes this connection explicit: the agent is directly penalised for poor predictive fitness at every moment, not merely for failing to improve.

---

**Algorithm 1** Curriculum Selection via Normalised Discounted Soft Q-Learning

---

1: **Initialise:** LM parameters  $\theta$ ; Q-network parameters  $\phi$ ; target parameters  $\bar{\phi} \leftarrow \phi$

2: // Bootstrap step (no  $Q$  update)

3:  $C_0 \leftarrow \text{NEXTBATCH}(\text{stream}, N)$

4:  $H_0 \leftarrow \text{DETACHEDHIDDEN}(\theta, C_0)$   $\triangleright [N, d] — \text{no grad into } \theta$

5:  $q_0 \leftarrow Q_\phi(H_0)$   $\triangleright [N, 1]$

6:  $\pi_0 \leftarrow \text{SOFTMAX}(q_0/\beta)$   $\triangleright \text{Boltzmann policy weights}$

7:  $a_0 \leftarrow \text{SAMPLE}(\pi_0)$

8:  $h_{\text{prev}} \leftarrow H_0[a_0]$   $\triangleright \text{Store hidden state of selected action}$

9:  $\theta \leftarrow \theta - \alpha \nabla_\theta \sum_i \pi_{0,i} \ell(\theta, C_0[i])$   $\triangleright \text{Mixture LM update, clipped to } \|\nabla\| \leq G$

10:  $r_{\text{prev}} \leftarrow \frac{1}{M} \log p_\theta(\hat{D}_0)$   $\triangleright \text{Reward after update}$

11: **for**  $k = 1, 2, \dots$  **do** // Main loop

12:    $C_k \leftarrow \text{NEXTBATCH}(\text{stream}, N)$

13:    $H_k \leftarrow \text{DETACHEDHIDDEN}(\theta, C_k)$   $\triangleright \text{Serves dual purpose below}$

14:    $q_k \leftarrow Q_\phi(H_k)$

15:   //  $Q$  update for previous transition (target network for bootstrap)

16:    $\bar{q}_k \leftarrow Q_{\bar{\phi}}(H_k)$   $\triangleright \text{Target network Q-values, no gradient}$

17:    $V_k \leftarrow \beta [\log \sum_{a'} \exp(\bar{q}_k[a']/\beta) - \log N]$   $\triangleright \text{Normalised soft value from target network}$

18:    $y \leftarrow (1 - \gamma) r_{\text{prev}} + \gamma V_k$   $\triangleright \text{Normalised TD target, stop-gradient}$

19:    $\hat{q} \leftarrow Q_\phi(h_{\text{prev}})$   $\triangleright \text{Re-evaluate previous action with current } \phi$

20:   Update  $\phi$  to minimise  $\frac{1}{2}(\hat{q} - y)^2$ , clipped to  $\|\nabla\| \leq G$

21:    $\bar{\phi} \leftarrow \tau \phi + (1 - \tau) \bar{\phi}$   $\triangleright \text{Polyak update}$

22:   // Action selection for current step (for  $Q$ -learning backup)

23:    $\pi_k \leftarrow \text{SOFTMAX}(q_k/\beta)$   $\triangleright \text{Boltzmann policy weights}$

24:    $a_k \leftarrow \text{SAMPLE}(\pi_k)$

25:    $h_{\text{prev}} \leftarrow H_k[a_k]$

26:   // LM mixture update and reward

27:    $\theta \leftarrow \theta - \alpha \nabla_\theta \sum_i \pi_{k,i} \ell(\theta, C_k[i])$   $\triangleright \text{Mixture LM update, clipped to } \|\nabla\| \leq G$

28:    $r_k \leftarrow \frac{1}{M} \log p_\theta(\hat{D}_k)$   $\triangleright \text{Fresh } \hat{D}_k \subset D, \text{ reward after update}$

29:    $r_{\text{prev}} \leftarrow r_k$

30: **end for**

---