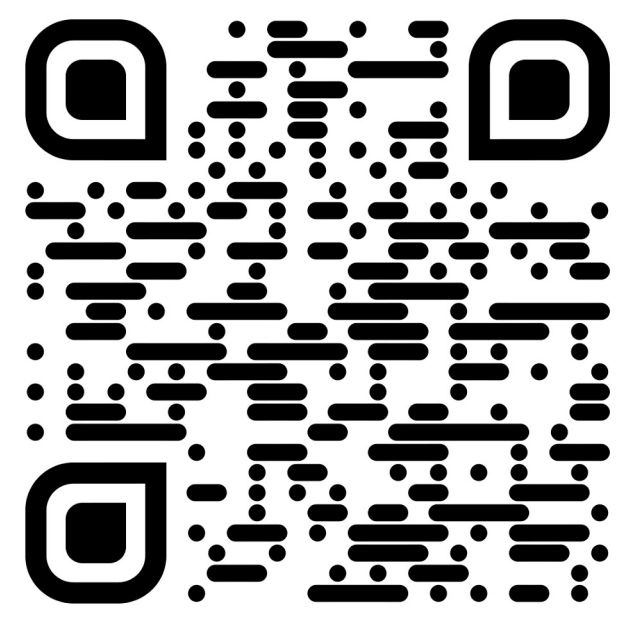
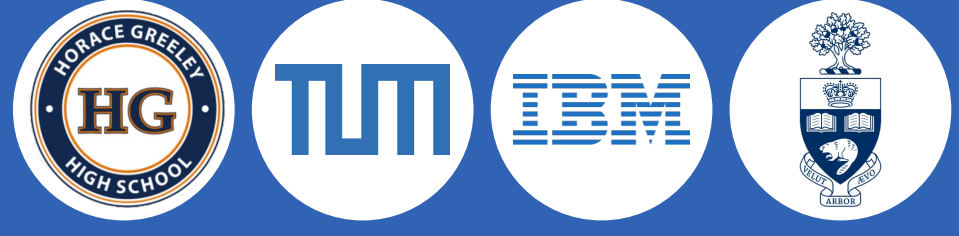


MESS+: Energy-Optimal Inferencing in Language Model Zoos with Service Level Guarantees



Ryan Zhang¹, Herbert Woisetschlager², Shiqiang Wang³, Hans-Arno Jacobsen⁴
 1 - Horace Greeley High School, 2 - Technical University of Munich, 3 - IBM Research, 4 - University of Toronto



Research Question: Can we select appropriate models from the model zoo to ensure energy efficiency while satisfying service level agreements (SLAs)?

1 Introduction

- Deep learning infrastructure providers and end users are confronted with an abundance of models (**model zoo**) for language modeling tasks.
- Related to two major research directions: **Dynamic Inference** and **Inference Request Scheduling**.
- Our work, MESS+, automatically selects a readily pre-trained model.
- Since MESS+ routes inference requests to different models, it can build on top of existing scheduling techniques.

We are solving a Tri-fold Problem

- End-users primarily care about *correct model output*
- Inference endpoint providers prioritize *low operating costs*
- Enterprise use-cases require consistent high quality model output while keeping costs in check through **Service Level Agreements (SLAs)**

2 Methodology



Algorithm 1: Selecting the Model with Energy-optimal Service level Guarantees (MESS+)

Input: $T; V; \alpha; c; \{E_m(t) : \forall m, t\}$; learning rate $\eta > 0$
Output: Outputs of models chosen for all t

- Initialize $Q(1) \leftarrow 0$; predictor parameters \mathbf{x}_m to a common random vector for all m ; $k \leftarrow 1$;
- for $t \leftarrow 1$ to T do
- Compute $p_t \leftarrow \min\left(1, \frac{c}{\sqrt[3]{t}}\right)$; **Exploration probability over time with cube root decay**
- Sample $\mathcal{X}_t \sim \text{Bernoulli}(p_t)$;
- if $\mathcal{X}_t = 1$ then
 - // Exploration
 - foreach $m \in \{1, 2, \dots, M\}$ do
 - Obtain true accuracy $A_m(t)$;
 - $\mathbf{x}_{m,t+1} \leftarrow \mathbf{x}_{m,t} - \eta \nabla_{\mathbf{x}} (\hat{A}(\mathbf{x}_{m,t}, \mathbf{a}_t) - A_m(t))^2$;
- $m^* \leftarrow \arg \max_m A_m(t)$;
- else **Decision problem for each request**
- $m^* \leftarrow \arg \min_m V \cdot E_m(t) + Q(t) \cdot (\alpha - \hat{A}_m(t))$; **A virtual queue is used to capture SLA violations**
- $\mathbf{x}_{m,t+1} \leftarrow \mathbf{x}_{m,t}$;
- Get output from model m^* and its accuracy $A_{m^*}(t)$;
- // Virtual queue update
- $Q(t+1) \leftarrow \max\{0, Q(t) + \alpha - A_{m^*}(t)\}$;

Overall Control Problem:

- E Energy consumption (joules) for a model.
- A Accuracy of a model's response (requires feedback signal).
- y Binary variable; 1 if a model is chosen.
- α SLA minimum accuracy requirement.

$$\min_{\{y_m(t) : \forall t, m\}} \frac{1}{T} \sum_{t=1}^T \sum_{m=1}^M y_m(t) E_m(t),$$

$$\text{s.t. } \frac{1}{T} \sum_{t=1}^T \sum_{m=1}^M y_m(t) A_m(t) \geq \alpha,$$

$$\sum_{m=1}^M y_m(t) = 1, \forall t \in \{1, \dots, T\},$$

$$y_m(t) \in \{0, 1\}, \forall t, m,$$

Challenges:

- Objective and constraints are correlated over time
- Characteristics of future requests cannot be predicted

3 Accuracy Predictor

MSE Objective for accuracy predictor:

$$L(\mathbf{x}_m) = \mathbb{E}_{\mathbf{a}_t} \left(\hat{A}(\mathbf{x}_m, \mathbf{a}_t) - A_m(t) \right)^2$$

To predict the accuracy, we train a predictor for K steps. The convergence upper bound is:

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla L(\mathbf{x}_{m,t_k})\|^2 \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$$

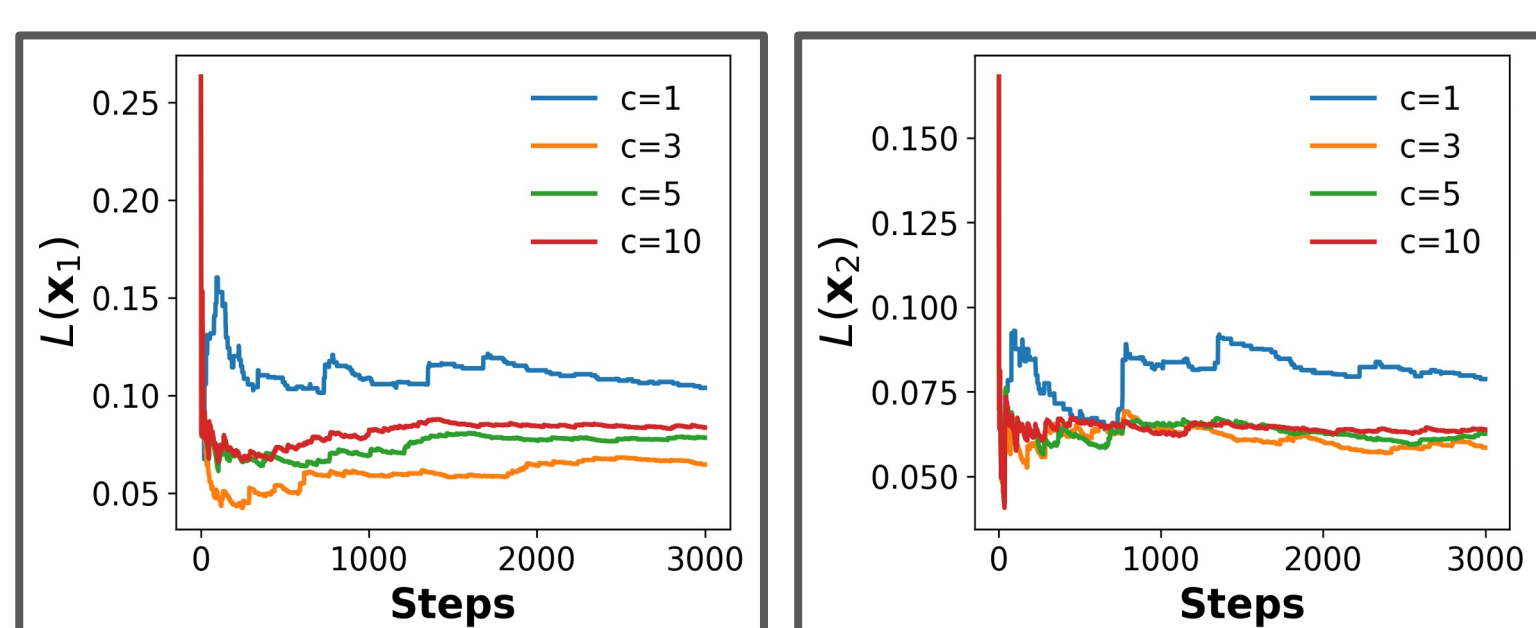
SGD convergence bound for our choice of p_t :

$$\mathbb{E}[K] = \Theta(T^{\frac{2}{3}}) \quad \mathcal{O}\left(1/\sqrt[3]{T}\right)$$

Upper bound of avg. add'l energy consumption: $\frac{1}{T} \sum_{k=1}^K E = \frac{KE}{T}$

For our choice of p_t : $\mathcal{O}\left(E/\sqrt[3]{T}\right)$

Our linear classifiers rapidly learn to predict their corresponding models' accuracy



4 MESS+ ensures energy minimal SLA compliance

- MESS+ satisfies the SLA with requirement α while consuming the least energy among all compliant strategies
- Analyzing V reveals that its optimal value is inversely related to the minimum service requirement α
- Balancing level of exploration is important as large c implies more training of \mathbf{x}_m but can also lead to overfitting on incoming requests

MESS+ reduces the energy consumption of an inference service by up to 2.5x

Model	WMT14 ($\alpha = 0.52$)			CNNDailyMail ($\alpha = 0.315$)		
	Accuracy (BLEU)	Energy (Joules)	Meets α	Accuracy (ROUGE1)	Energy (Joules)	Meets α
TinyLlama	49.1 \pm 0.6	44.639 \pm 0.6	No	30.9 \pm 0.3	142.080 \pm 1.4	No
Llama-2 13B	55.1 \pm 0.4	527.870 \pm 5.1	Yes	32.2 \pm 0.3	750.285 \pm 7.5	Yes
Random with constraint	52.0 \pm 0.0	280.426 \pm 3.0	Yes	31.5 \pm 0.0	416.466 \pm 4.3	Yes
MESS+ ($V = 0.1, c = 3$)	52.2 \pm 0.2	149.399 \pm 1.3	Yes	31.5 \pm 0.1	163.836 \pm 1.5	Yes

