

Consistent Accelerated Inference via Confident Adaptive Transformers

Tal Schuster*, Adam Fisch*, Tommi Jaakkola, and Regina Barzilay (*= equal contribution)



Overview

Goal: Reduce the computational effort of multi-layered deep models while ensuring consistency with the original model \mathcal{F} that uses l layers.

Approach: A new model \mathcal{G} that can exit at layer $\leq l$ and guarantees:

$$\mathbb{P}(\mathcal{G}(X_{n+1}) = \mathcal{F}(X_{n+1})) \geq 1 - \epsilon$$

Specifically, our model is defined as:

$$\mathcal{G}(x; \tau) := \begin{cases} \mathcal{F}_1(x) & \text{if } \mathcal{M}_1(x) > \tau_1, \\ \mathcal{F}_2(x) & \text{else if } \mathcal{M}_2(x) > \tau_2, \\ \vdots & \\ \mathcal{F}_l(x) & \text{otherwise,} \end{cases}$$

Where \mathcal{M} are confidence scores and τ are stopping thresholds.

Challenges:

- 1) *What is a good confidence score?* Meta predictor
- 2) *How to calibrate the thresholds?* Conformal calibration over inconsistent layers

Conformalized Early Exit

We look at **inconsistent** layers:

$$\mathcal{I}(x) := \{i : \mathcal{F}_i(x) \neq \mathcal{F}(x)\}, \quad i \in [1, l-1]$$

and obtain a conservative prediction set \mathcal{C}_ϵ s.t.:

$$\mathbb{P}(\mathcal{I}(X_{n+1}) \subseteq \mathcal{C}_\epsilon(X_{n+1})) \geq 1 - \epsilon.$$

The first layer in the complement set provides:

$$\mathbb{P}(\mathcal{F}_K(X_{n+1}) = \mathcal{F}(X_{n+1})) \geq 1 - \epsilon,$$

where $K := \min \{j : j \in \mathcal{C}_\epsilon^c(X_{n+1})\}$.

Independent calibration: For each layer, compute the empirical distribution over a calibration set:

$$v_k^{(1:n,\infty)} := \{\mathcal{M}_k(x_i) : x_i \in \mathcal{D}_{\text{cal}}, \mathcal{F}_k(x_i) \neq \mathcal{F}(x_i)\} \cup \{\infty\}.$$

Take the quantile after MHT correction:

$$\tau_k^{\text{ind}} = \text{Quantile}(1 - \alpha_k, v_k^{(1:n,\infty)}).$$

Shared Calibration: Calibrate for the *worst-case* across inconsistent layers (the maximum score):

$$m^{(1:n,\infty)} := \{\mathcal{M}_{\max}(x_i) : x_i \in \mathcal{D}_{\text{cal}}, \exists k \mathcal{F}_k(x_i) \neq \mathcal{F}(x_i)\} \cup \{\infty\};$$

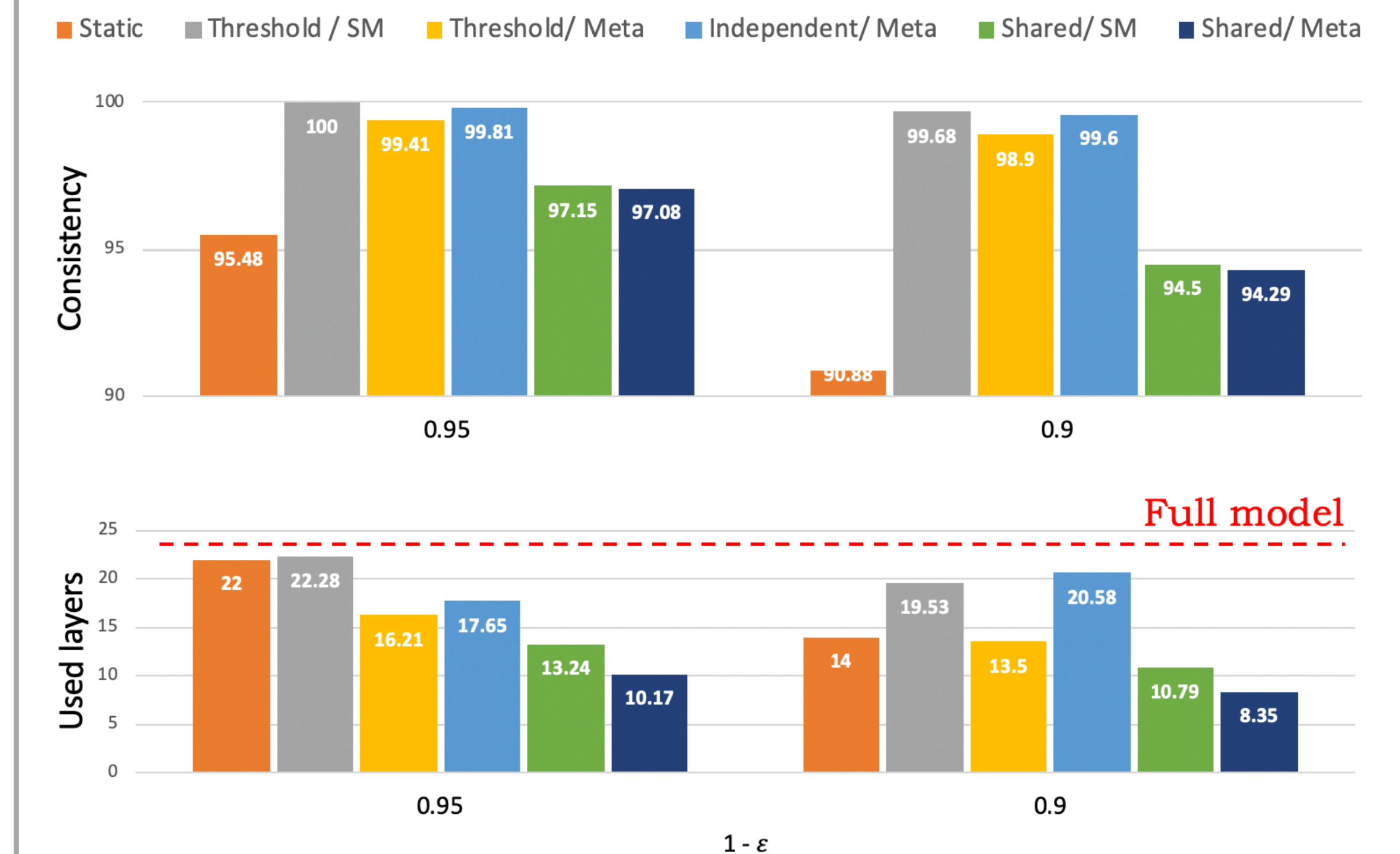
$$\tau^{\text{share}} = \text{Quantile}(1 - \epsilon, m^{(1:n,\infty)}).$$

Experiments

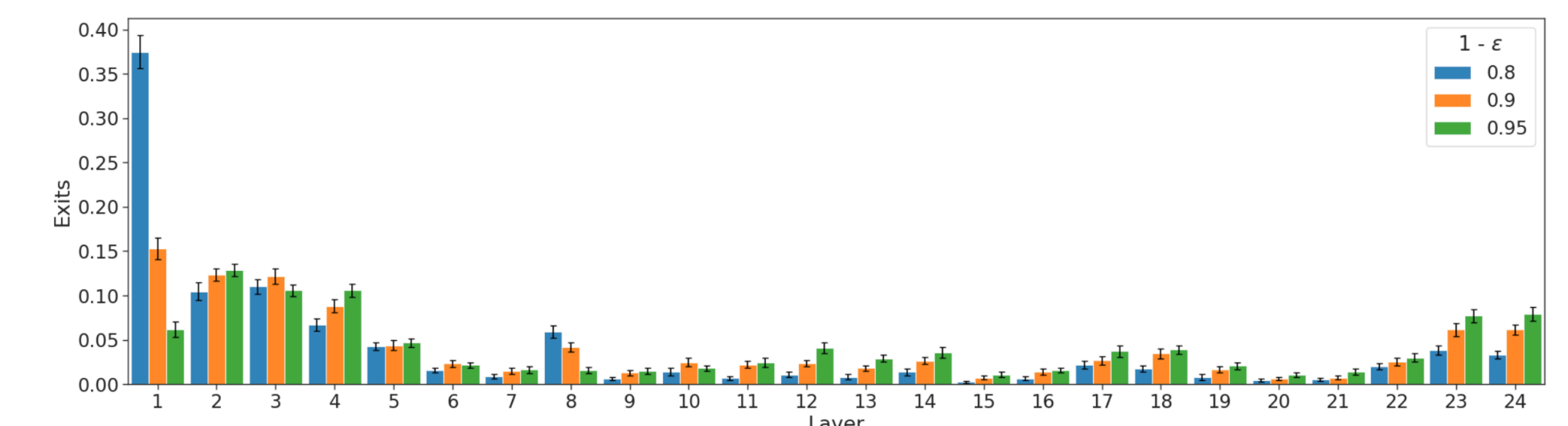
Experimental setting: we implement our Confident Adaptive Transformers (CATs) on top of popular models (e.g. Albert). We evaluate on four NLP tasks covering both classification and regression.

Evaluated on IMDB, VitaminC, AG news, STS-B

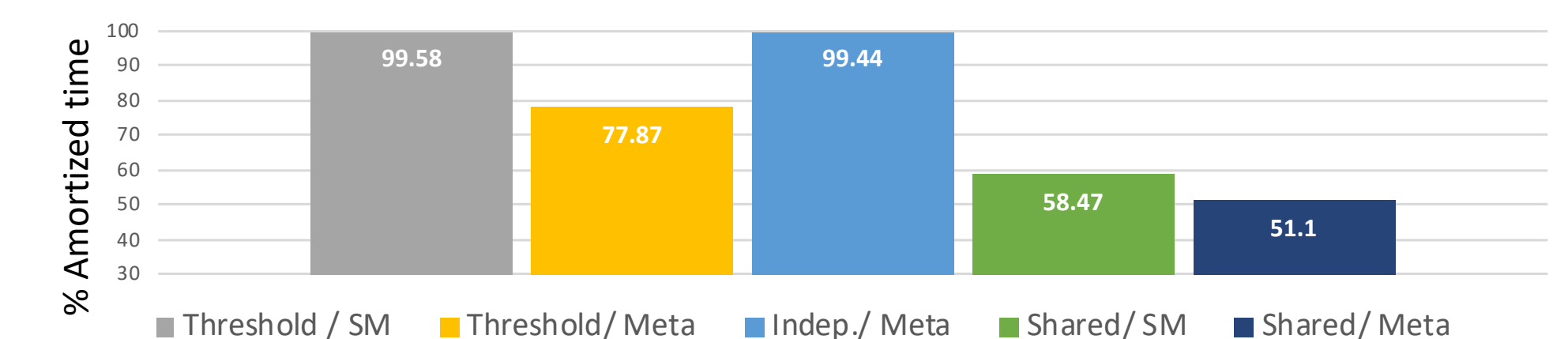
Results on AG news:



The distribution of exit layers is dynamically controlled by the user-defined tolerance level:



Amortized time ($1 - \epsilon = 0.9$):



Meta confidence per Transformer layer

