

MD-DiT: Step-aware Mixture-of-Depths for Efficient Diffusion Transformers

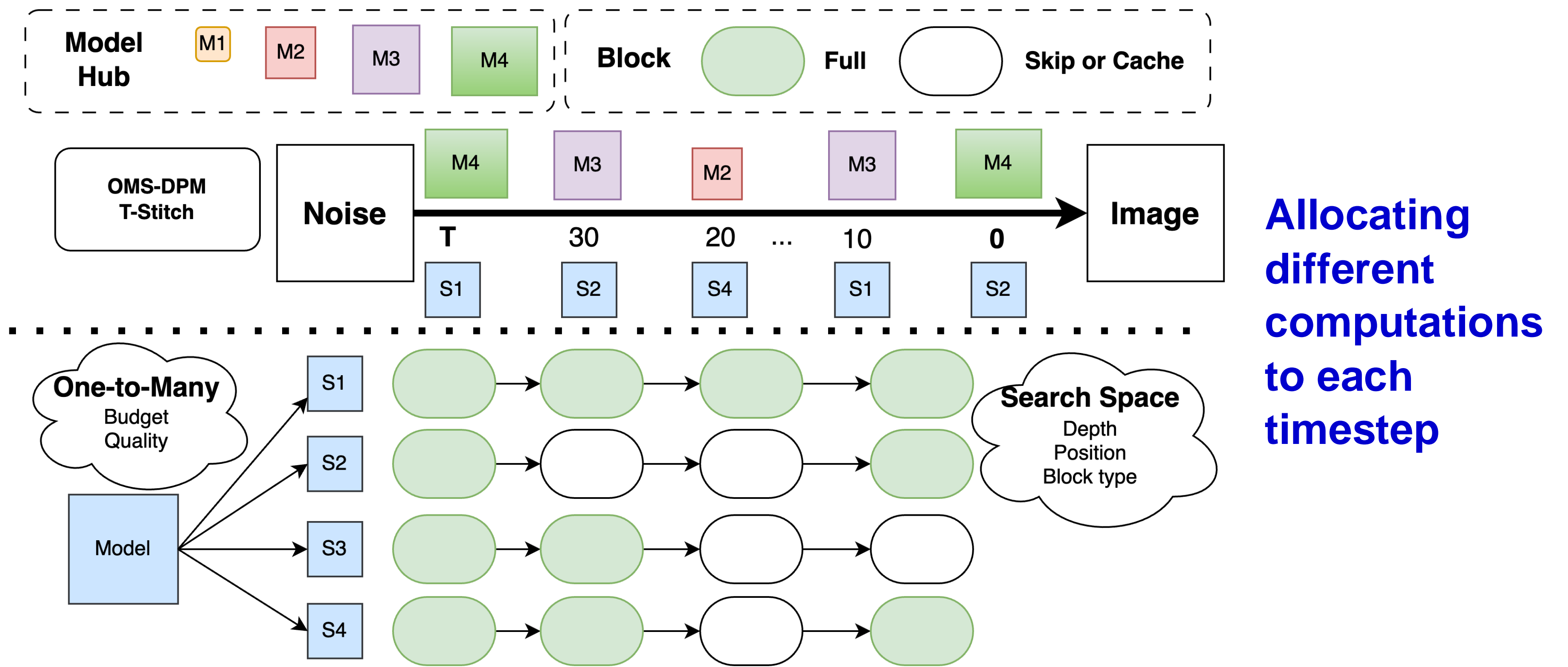
Mingzhu Shen^{1*}, Pengtao Chen^{2*}, Peng Ye^{3,4}, Guoxuan Xia¹, Tao Chen², Christos-Savvas Bouganis¹, Yiren Zhao¹
¹Imperial College London, ²Fudan University, ³CUHK, ⁴Shanghai AI Lab, *Equal Contribution

Framework

MD-DiT, a **one-to-many** unified framework that realizes a **mixture-of-depths across different timesteps** via the incorporation of block **skipping and caching** techniques.

Our contributions are as follows:

- We introduce MD-DiT, a framework that combines block skipping and caching to create a **mixture-of-depths across timesteps**.
- We explore **depth allocation strategies** for each timestep and use a **gradient-free search method** to identify a more compact model, improving diffusion transformer acceleration.



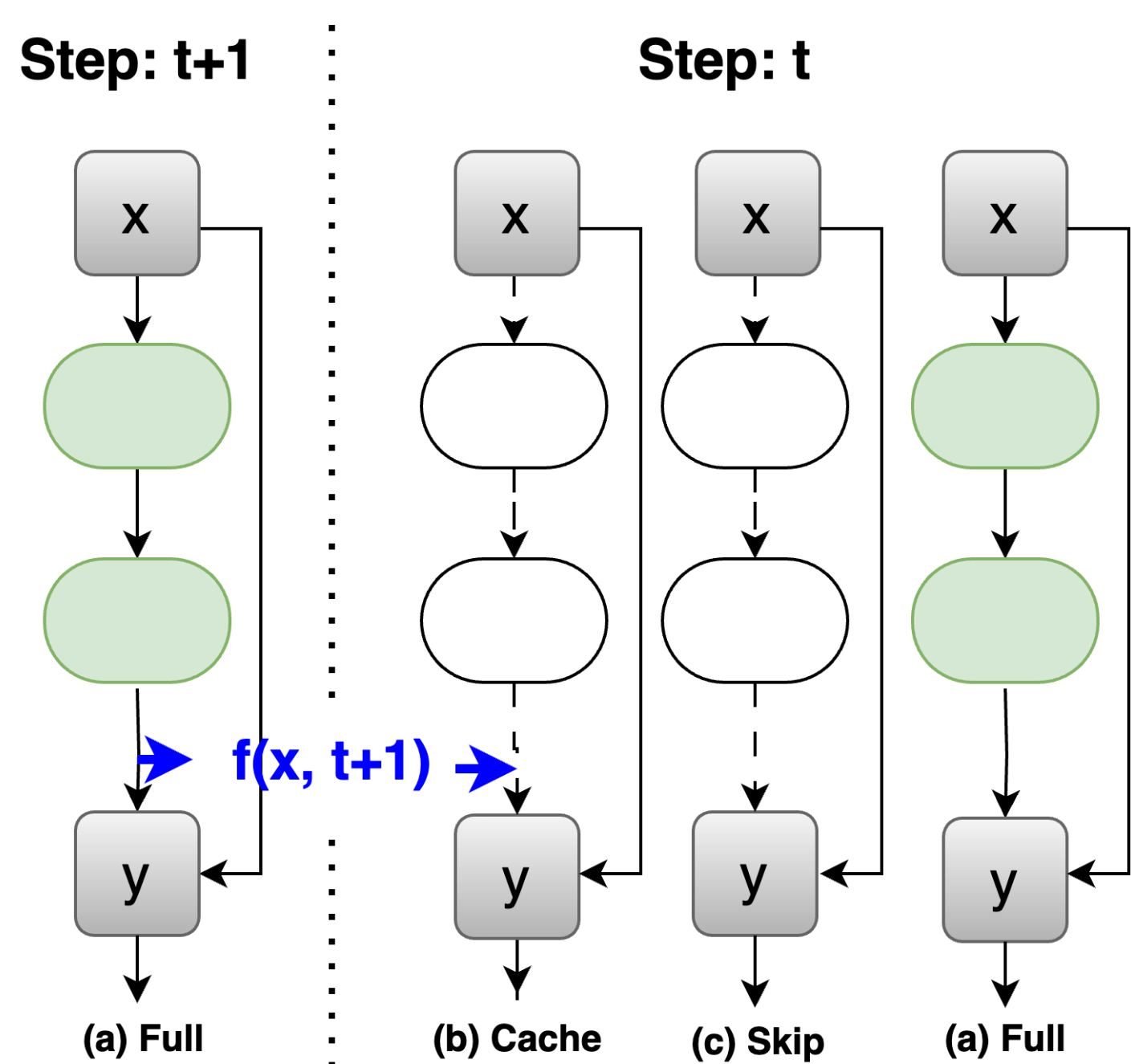
Search Space

Search Insights

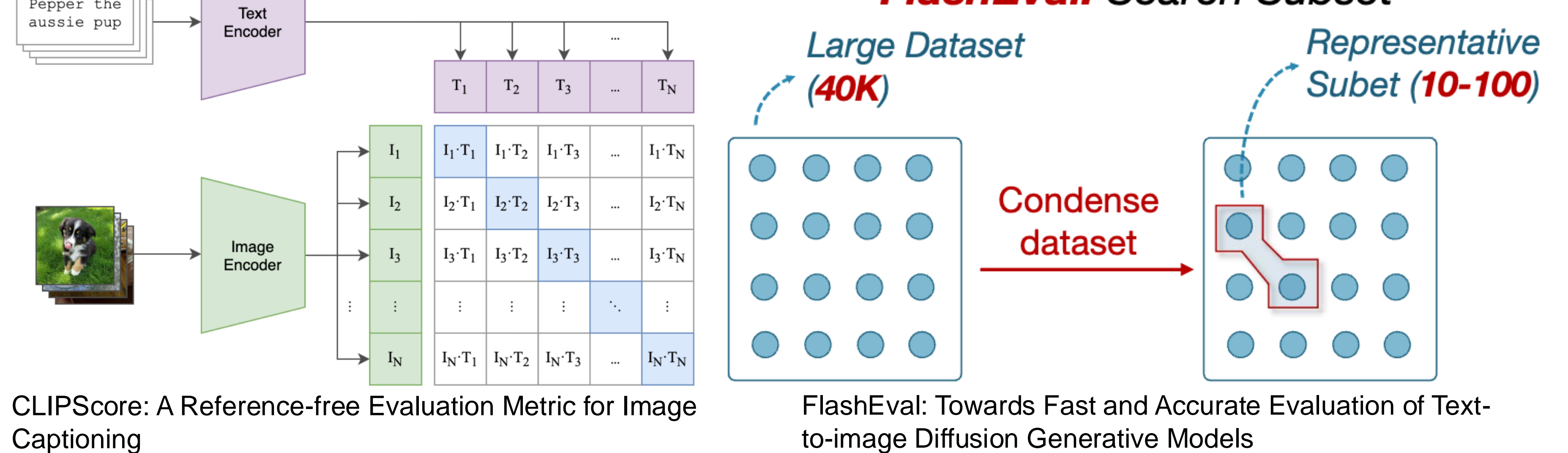
Search to Skip or Cache or Compute

How to Search Efficiently: **Clip Score** on a small **Search SubSet**

Gradient-Free Search: Covariance Matrix Adaptive Evolution Strategy



FlashEval: Search Subset



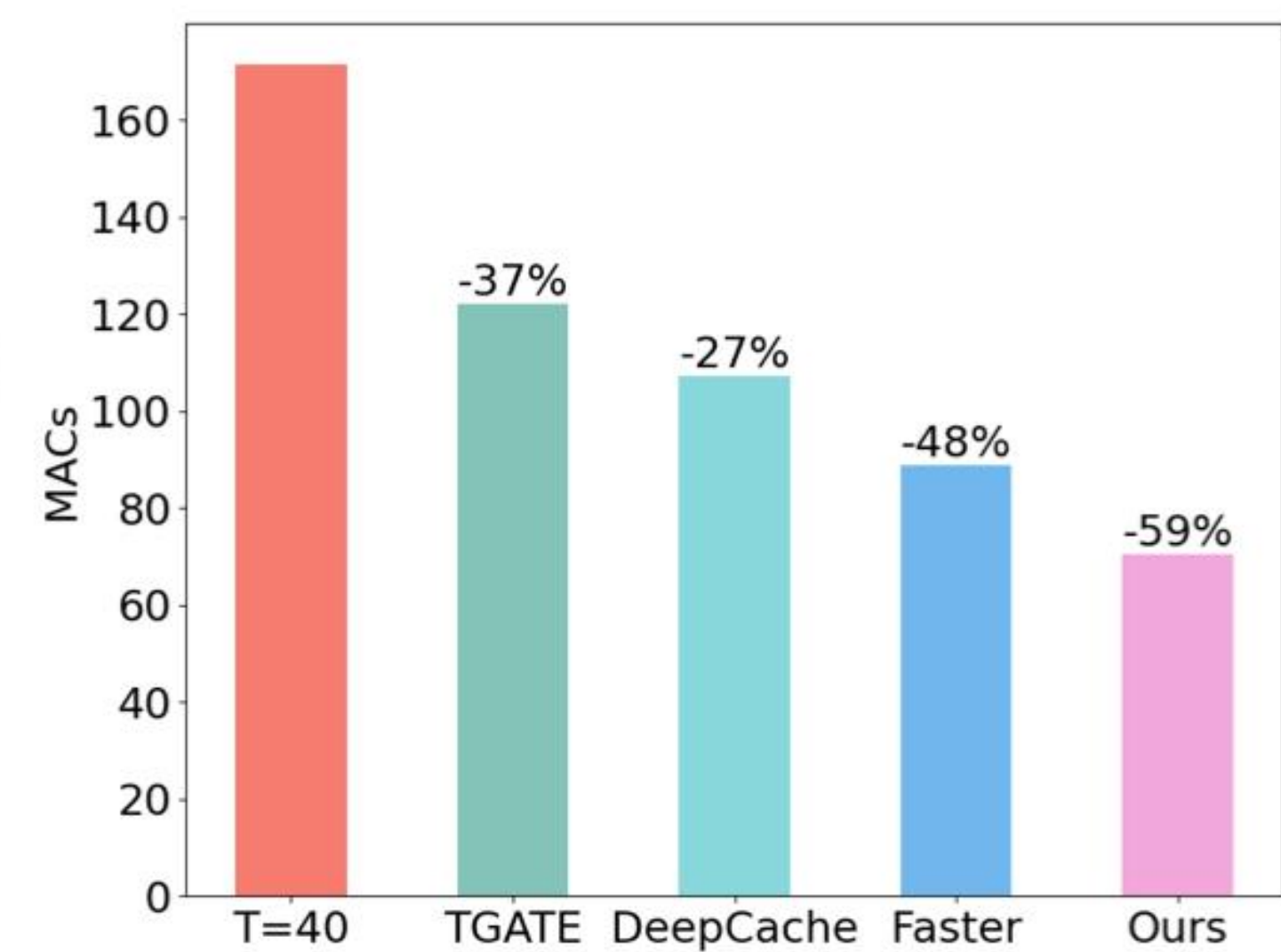
Experiments

Through extensive experiments, we have successfully compressed the **LCM-4Step** model with a **20% reduction** in Multiple-Accumulate Operations (MACs). This achievement is further amplified in a **40-step** setting, where we have accomplished a **59% reduction**.

(a) LCM Pixart-Alpha for MSCOCO-2017



(b) DiT-XL on ImageNet



(c) Pixart-Alpha for MSCOCO-2017 with T=40. The correspondent Clip Score for each method is 30.45, 29.9, 30.2, 30.4, 30.4. respectively.