

# Grape: Practical and Efficient Graph-based Executions for Dynamic Deep Neural Networks on GPUs



Bojian Zheng<sup>1,2,3</sup>, Cody Hao Yu<sup>4</sup>, Jie Wang<sup>4</sup>, Yaoyao Ding<sup>1,2,3</sup>, Yizhi Liu<sup>4</sup>, Yida Wang<sup>4</sup>, Gennady Pekhimenko<sup>1,2,3</sup>

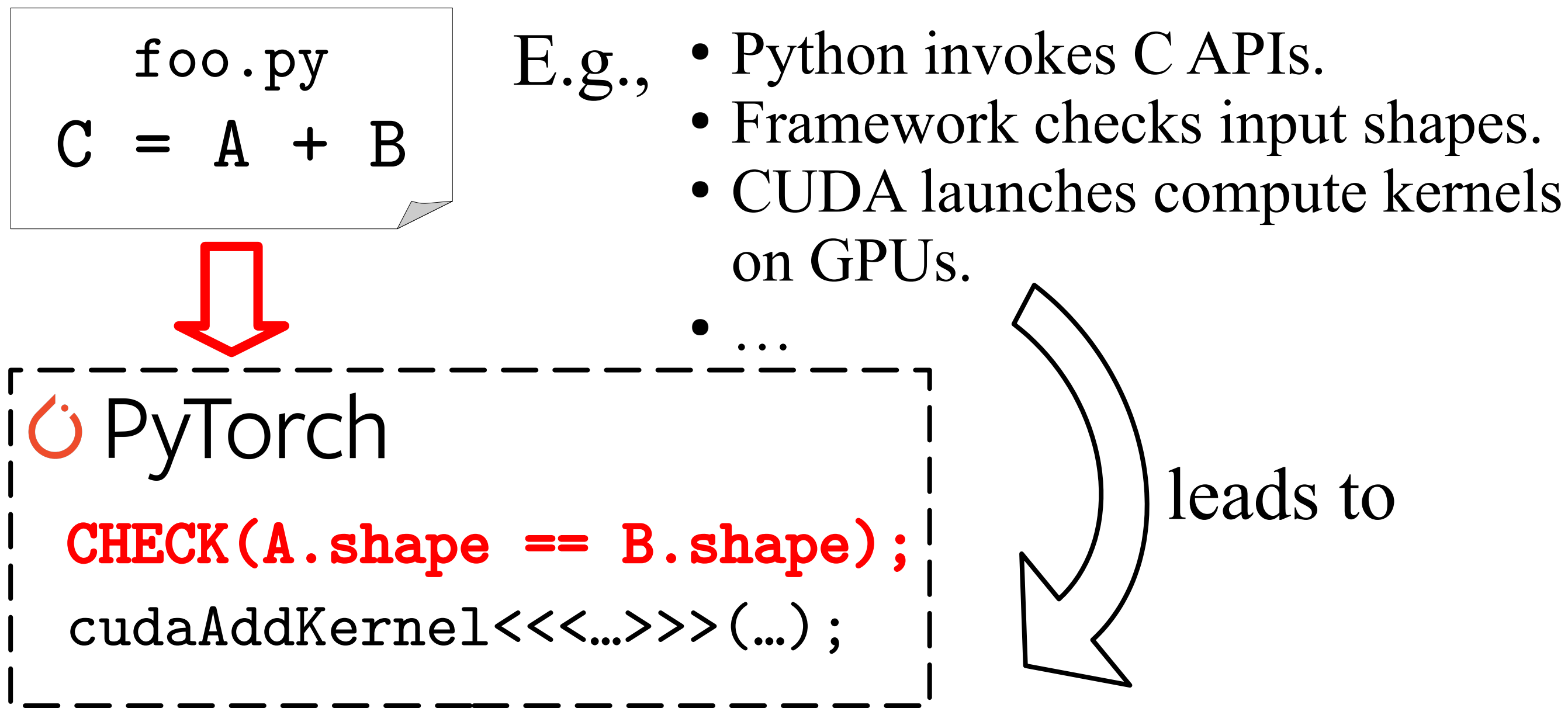
<sup>1</sup>CentML <sup>2</sup>University of Toronto <sup>3</sup>Vector Institute <sup>4</sup>AWS

<https://github.com/UofT-EcoSystem/Grape-MICRO56-Artifact>

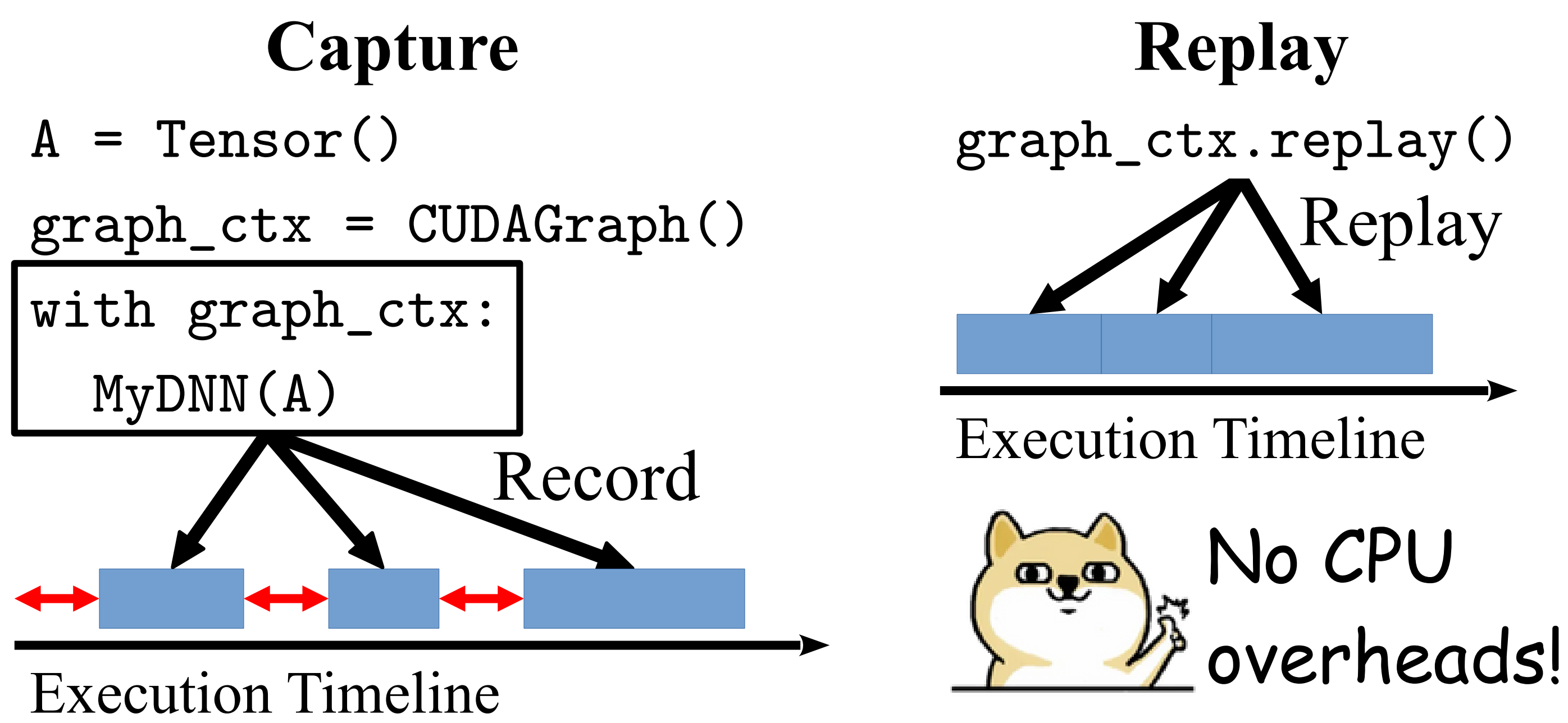


## 1. Background: CUDA Graphs

Ubiquitous CPU overheads in machine learning systems:

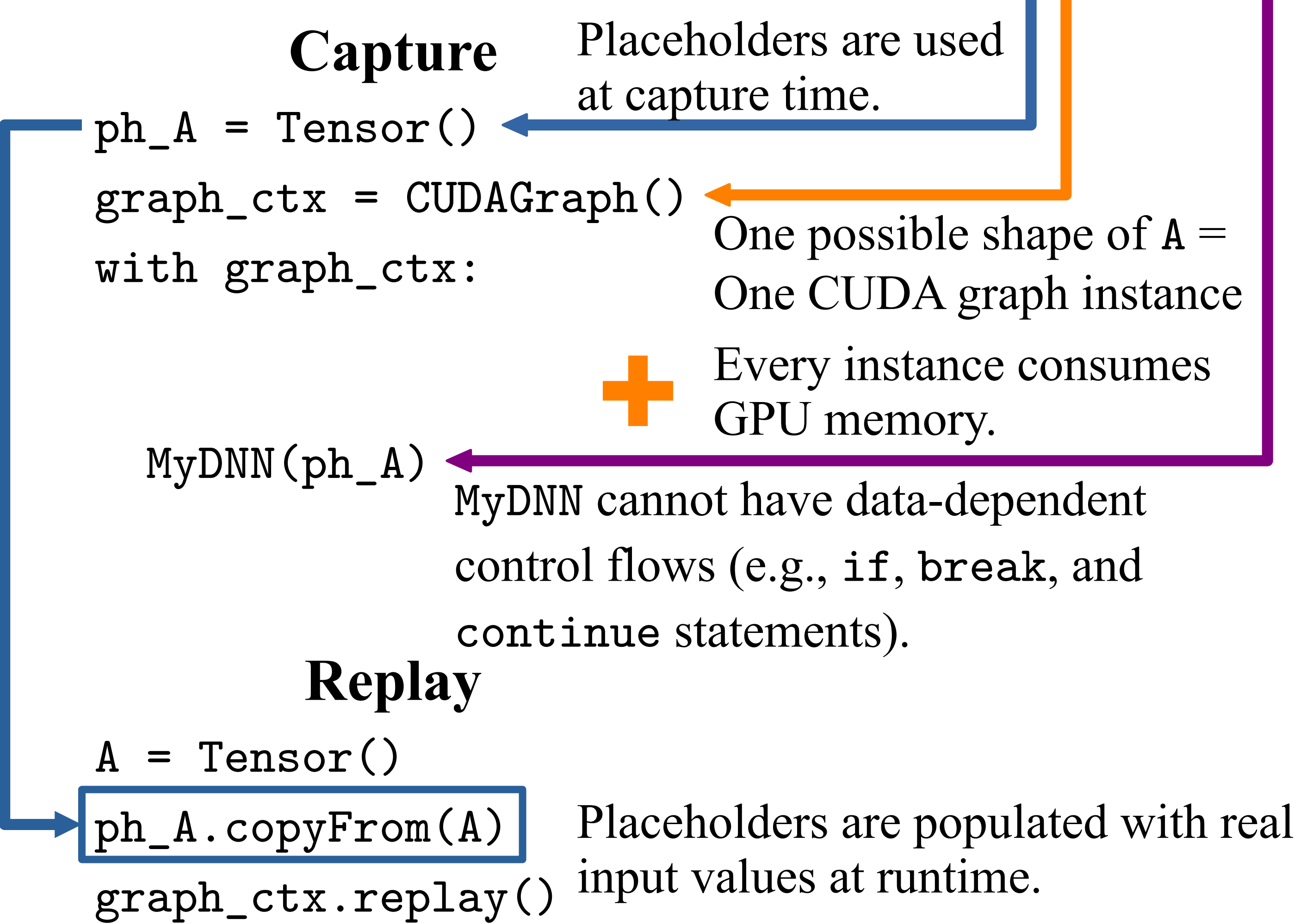


CUDA graphs remove CPU overheads by capturing effective GPU operations and then replaying them:



## 2. Challenges posed by CUDA Graphs

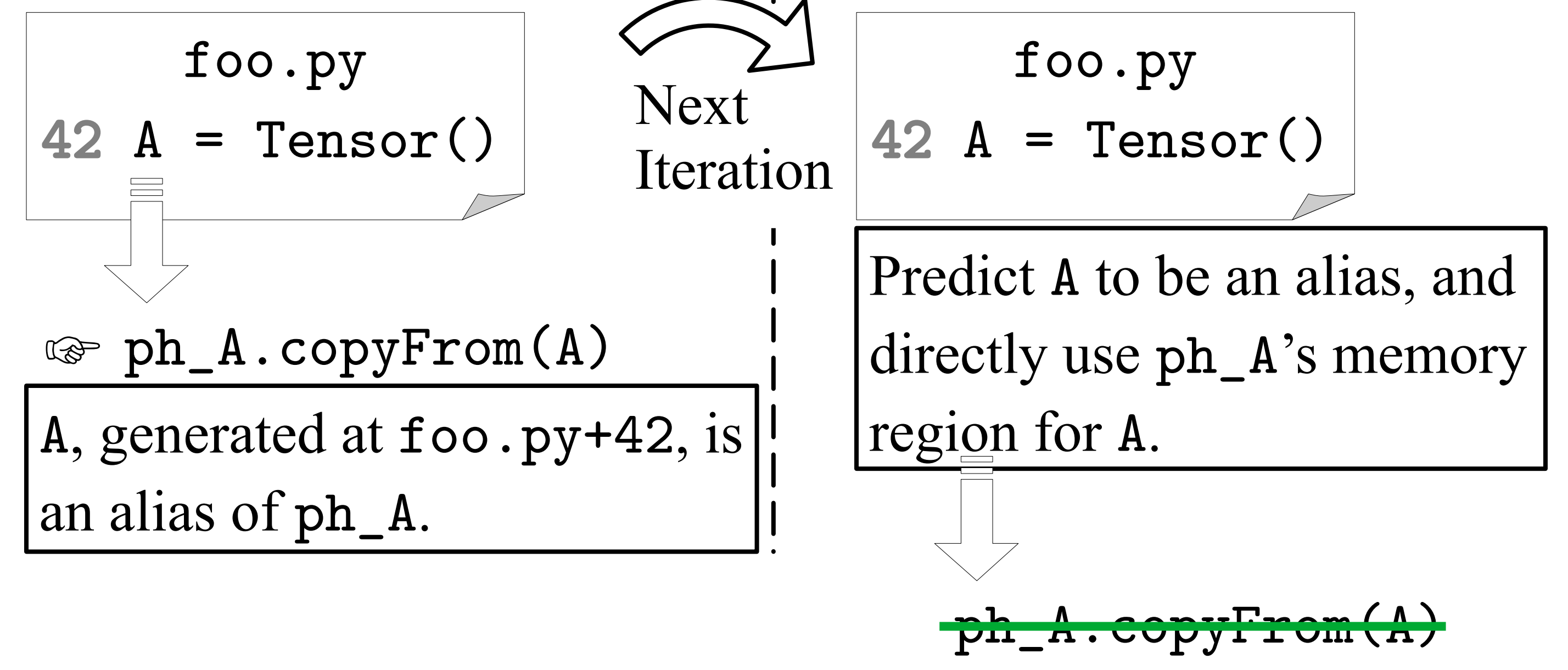
CUDA graphs request computations to be frozen.



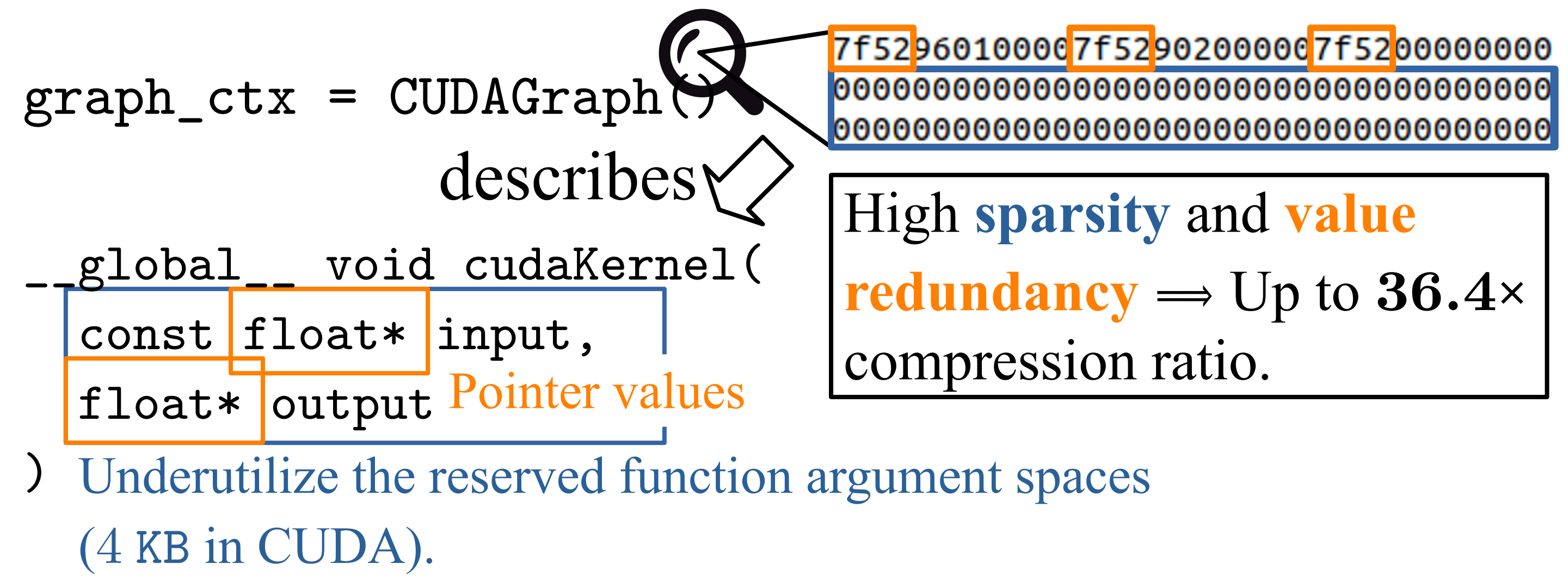
- 1 Extra data movements into placeholders incur runtime overheads (up to 13%).
- 2 Huge GPU memory consumption to efficiently support dynamic-shape workloads (20-100 GB).
- 3 No support for data-dependent control flows.

## 3. Grape's Key Ideas

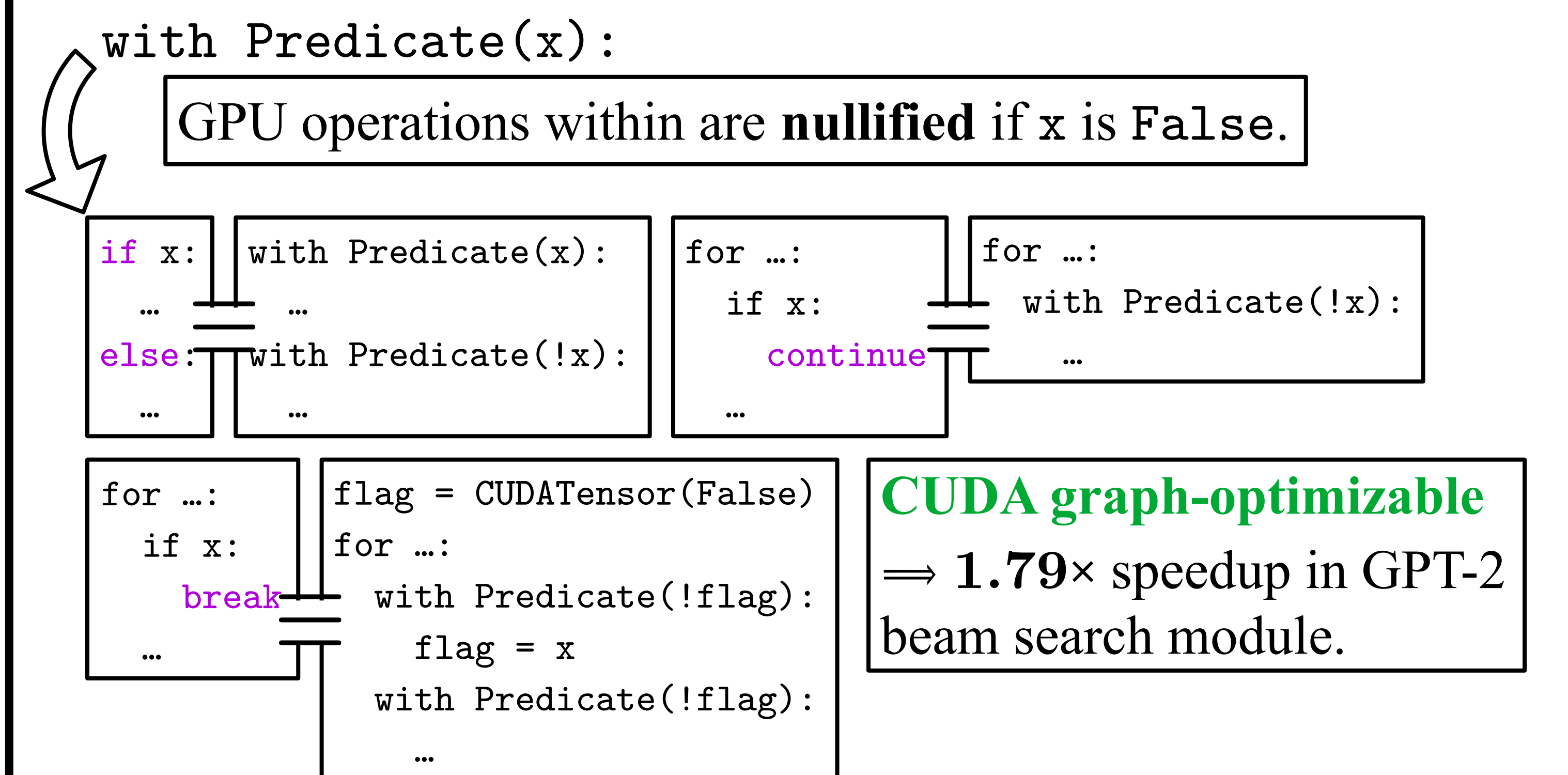
1 Alias Prediction



2 Metadata Compression



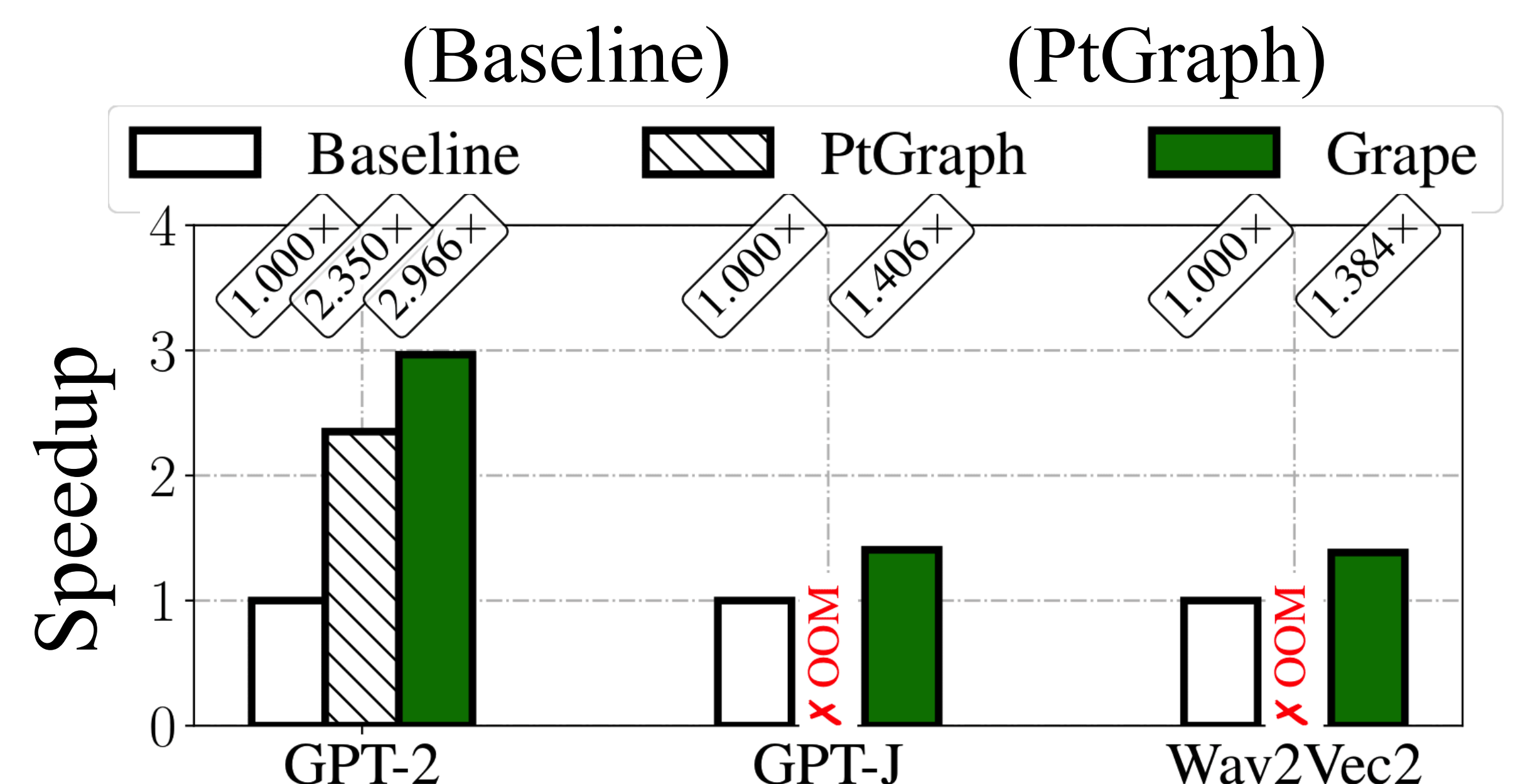
3 Predication Contexts



## 4. Evaluation

Infrastructure: NVIDIA RTX 3090 and A100 with PyTorch ver. 1.12 and CUDA ver. 12.0. Applications: GPT-2<sup>137M</sup>, GPT-J<sup>6B</sup>, and Wav2Vec2<sup>Base</sup> from Transformers.

Compared with PyTorch and Graphs + PyTorch.



- On GPT-2, 2.97x/1.26x better than Baseline/PtGraph.
- On GPT-J and Wav2Vec2, up to 1.41x better than Baseline while PtGraph goes OOM.