# Memory saving
# for neural network trainng

Xunyi Zhao

# 01

Introduction

*Inría*

- Rotor: Rematerializing Optimally with pyTORch
  - Efficient with sequential models

*Inria*

- Rotor: Rematerializing Optimally with pyTORch
  - Efficient with sequential models
- Weight Offloading
  - Non-optimal but efficient

*Inria*

- Rotor: Rematerializing Optimally with pyTORch
  - Efficient with sequential models
- Weight Offloading
  - Non-optimal but efficient
- Rotor+checkmate(Ongoing)
  - Combined ideas to be general and efficient

*Inría*

# 02

Weight offloading

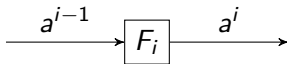Possibility of training large model in a single GPU

Only keep what is necessary

*Inria*

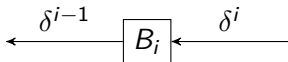Possibility of training large model in a single GPU

Only keep what is necessary

**Forward $i$**

- $a^i \leftarrow a^{i-1}, w_i$

$$\xrightarrow{a^{i-1}} \boxed{F_i} \xrightarrow{a^i}$$

## Backward $i$

- $\delta w_i \leftarrow a^{i-1}, \delta^i$;
- $\delta^{i-1} \leftarrow \delta^i, w_i$;
- $w_i \leftarrow w_i, \delta w_i$

$$\xleftarrow{\delta^{i-1}} \boxed{B_i} \xleftarrow{\delta^i}$$

*Inria*

Offload layers that cover most memory needs

## Greedy selection

- Assume infinite bandwidth
- Start: every layer in GPU
- Choose $w_i$ that covers the most memory overflow
- Iterate until memory fits
- A coefficient is introduced for single offloading

*Inria*

## Greedy schedule

- Insert prefetch and offload operations into $F$'s and $B$'s
- Data required earlier is prefetched first
- Data produced earlier is offloaded first
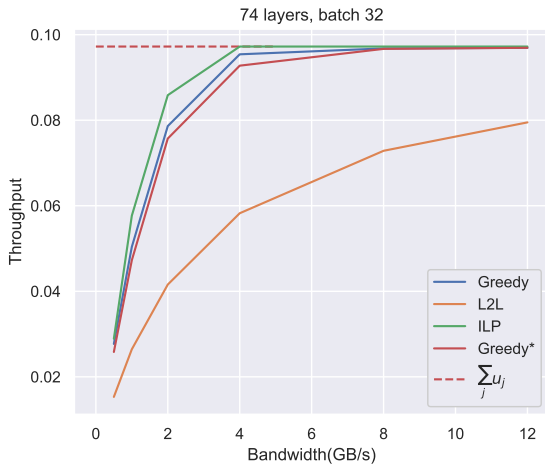- One particular case: $w_i$ is sent to CPU after $B_i$, but only deleted after $F_i$

*Inria*

Figure: Greedy* represents the Greedy algorithm while weights are not only offloaded once.

# 03

Rotor+checkmate

*Inria*

Checkmate provides an optimal solution for rematerializations.

## Application

- Can deal with any graphs
- ILP with $O(n^2)$ variables
- Implementations based on TensorFlow

*Inria*

Rotor is efficient, checkmate is general

Rotor is efficient, checkmate is general

**Combined ideas**

- One network can be divided into several blocks: e.g. GPT-2 → Decoder
- Each block is (probably) small enough for checkmate
- Rotor can solve the schedule based on the solutions from checkmate

*Inria*

- Checkmate on PyTorch: obtain the computation graph from any Torch Module

  - We can now open the box and generate the forward&backward graph

  - Each node in the graph can be executed by a torch function/python code

*Inria*