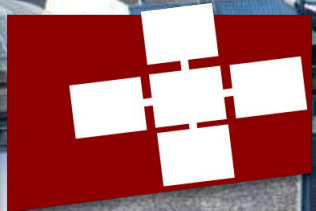


KAZUKI OSAWA^{1*}, SHIGANG LI², TORSTEN HOEFLER¹

1: ETH Zurich (*: currently with Google DeepMind), 2: Beijing University of Posts and Telecommunications

PipeFisher: Efficient Training of Large Language Models Using Pipelining and Fisher Information Matrices



Large Language Models

ChatGPT, <https://chat.openai.com/>

k What are Large Language Models?

Large Language Models (LLMs) are powerful artificial intelligence systems that are designed to understand and generate human-like text. These models, such as OpenAI's GPT-3, are trained on vast amounts of text data to learn the statistical patterns and structures of language.

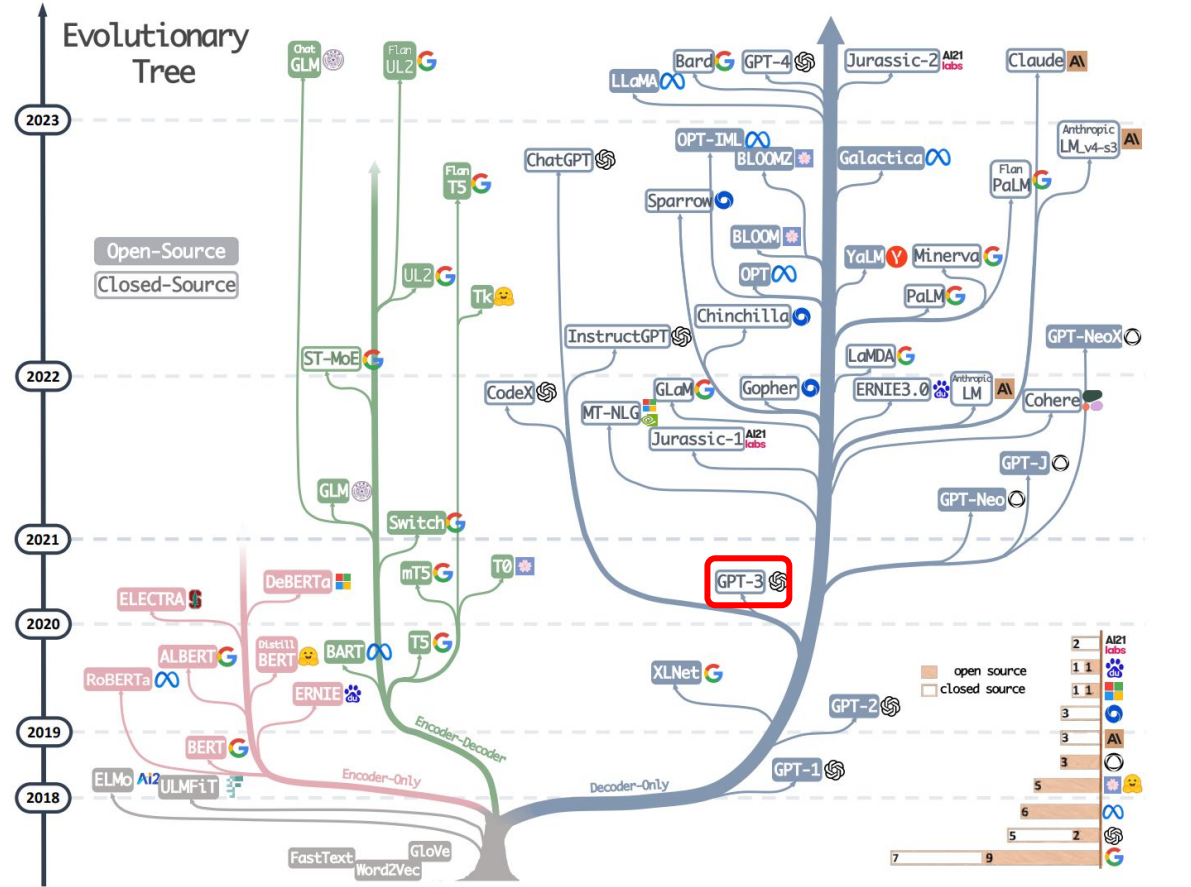
LLMs are based on deep learning techniques, particularly using a type of neural network called a transformer. Transformers allow LLMs to process and generate text in a hierarchical manner, taking into account the relationships between words, phrases, and sentences.

The training process for LLMs involves exposing the model to massive amounts of text from diverse sources such as books, websites, and articles. By predicting the next word in a sentence or completing a given prompt, the model learns to generate coherent and contextually appropriate responses. (GPT-3: 300B tokens)

LLMs have a wide range of applications. They can be used for natural language processing tasks such as text classification, sentiment analysis, language translation, summarization, and question answering. They can also be employed in creative tasks like generating human-like text, poetry, or even composing music.

GPT-3, for example, with 175 billion parameters, is one of the largest LLMs developed to date. These massive models have demonstrated impressive language understanding and generation capabilities, although they also raise concerns related to ethical considerations, biases, and potential misuse.

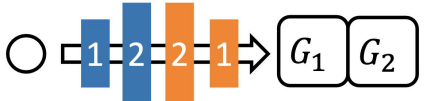
Image: Yang et al., "Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond", 2023.



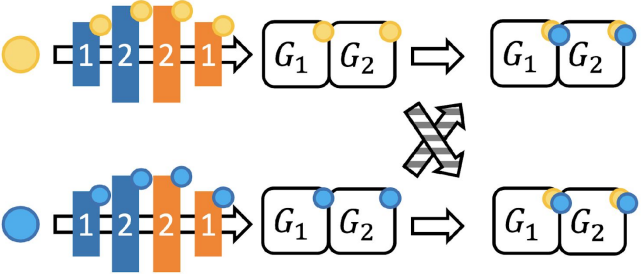
Data scale ↑ **> 1T tokens** + Model size ↑ **> 1T parameters** ⇨ Cost ↑ **> 40K A100 GPU days**
(e.g., Chinchilla, LLaMA) (e.g., GLaM) (e.g., BLOOM)

Parallelism for massive data and models

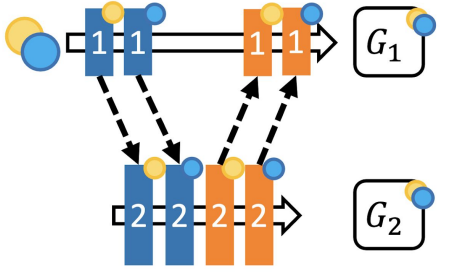
(i) No parallelism
(w/ 1 device)



(ii) Data parallelism
(w/ 2 devices)



(iii) Pipeline parallelism
(w/ 2 devices)



Forward pass

pre activation

$$h_l = W_l a_{l-1}$$

$$(a_0 = x)$$

activation

$$a_l = \phi(h_l)$$

Backward pass

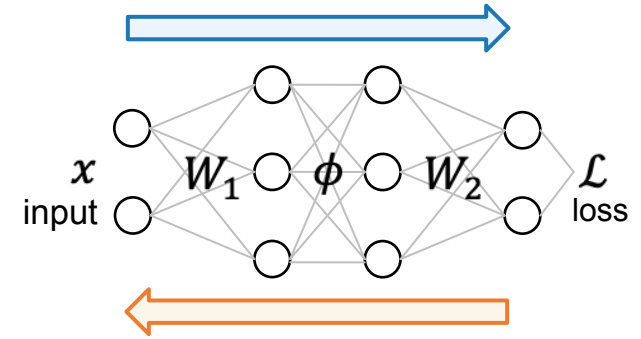
$$\frac{\partial \mathcal{L}}{\partial a_{l-1}} = W_l^T \frac{\partial \mathcal{L}}{\partial h_l}$$

error

$$e_l = \frac{\partial \mathcal{L}}{\partial h_l} = \phi' \left(\frac{\partial \mathcal{L}}{\partial a_l} \right)$$

gradient

$$G_l = \frac{\partial \mathcal{L}}{\partial W_l} = e_l a_{l-1}^T$$



- Mini- or micro-batch
- | | x x Forward/backward at x-th layer for ○
- q_x Quantity for x-th layer calculated for ○
- Collective communication
- - - → Point-to-point communication (send/recv)

and more model parallelism (ZeRO, Megatron, MoE, etc)

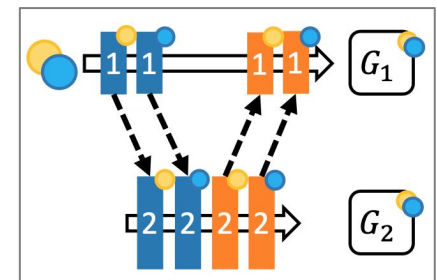
Bubbles in pipeline

- Pipelining creates **bubbles** of time in which accelerators become idle.
- The overhead of pipelining mainly comes from the **low utilization of accelerators**.
- 💡 We can assign **extra work** to the bubbles to gain **auxiliary benefits**.
- [Our approach] **PipeFisher** *automatically* assigns the **work of K-FAC** (a second-order optimization method based on the Fisher information matrix) to the bubbles for **accelerating training**.

(4 pipeline stages, 4 micro-batches)

(2 pipeline stages, 2 micro-batches)

(a) GPipe (Huang et al., 2018)



(b) PipeFisher for GPipe



Why second-order optimization?

First-order Optimization (gradient descent)

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)})$$

Second-order Optimization

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \mathbf{C}^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

Precondition the gradient by the **curvature matrix**

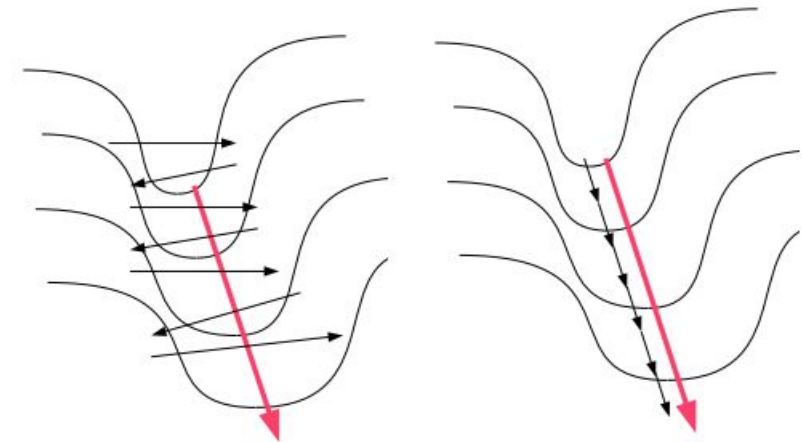
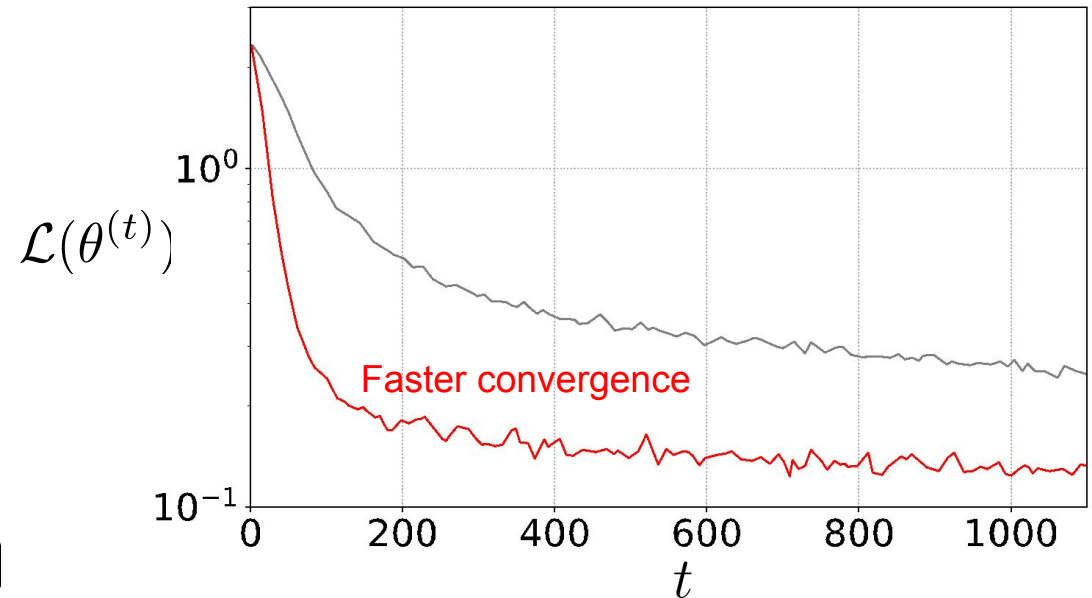


Figure from J. Martens, 2010

Why second-order optimization?

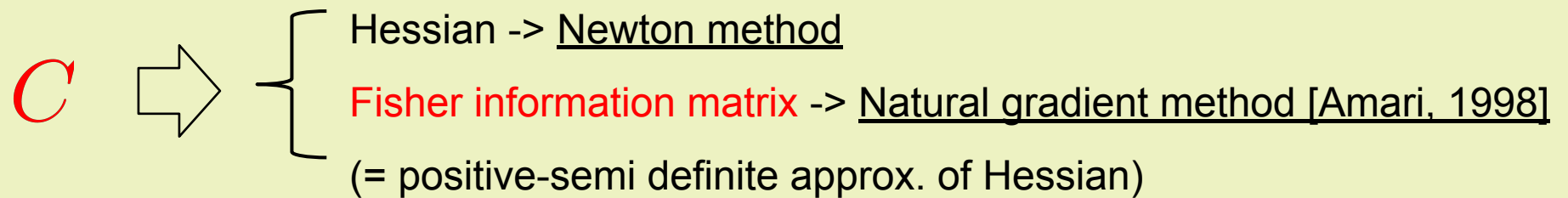
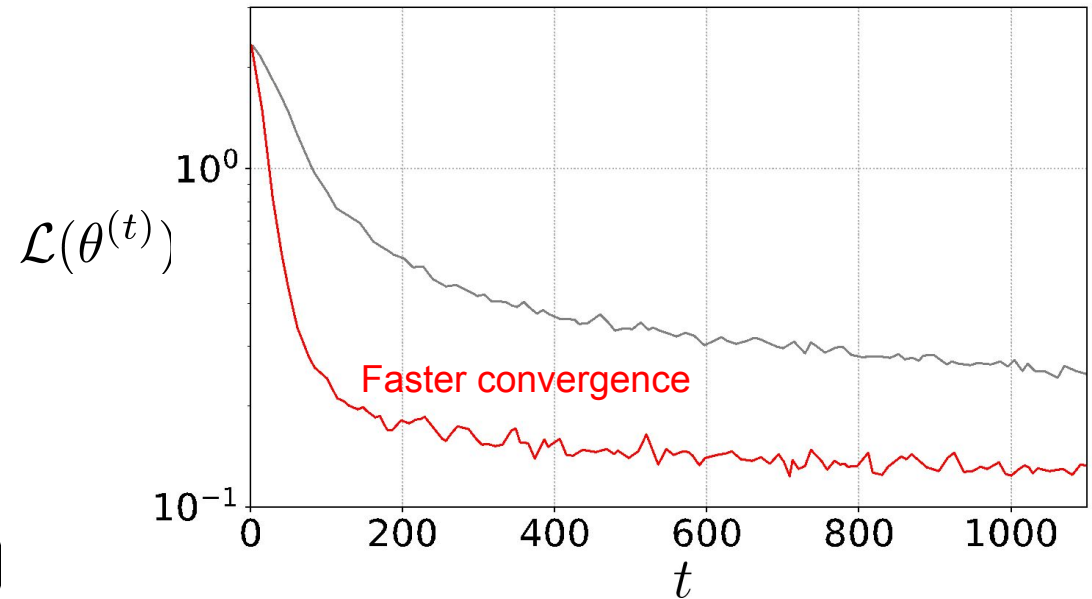
First-order Optimization (gradient descent)

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)})$$

Second-order Optimization

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \mathbf{C}^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

Precondition the gradient by the **curvature matrix**



Why is second-order optimization *unpopular* in DL?

First-order optimization

Second-order optimization

Forward/backward

$$\nabla \mathcal{L} \in \mathbb{R}^P$$

$$\nabla \mathcal{L} \in \mathbb{R}^P$$

Curvature

$$C \in \mathbb{R}^{P \times P}$$

Inverse

$$C^{-1} \in \mathbb{R}^{P \times P}$$

Precondition

$$C^{-1} \nabla \mathcal{L} \in \mathbb{R}^P$$

Overhead
 $\mathcal{O}(P^3)$

Update

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)})$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta C^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

Why is second-order optimization *unpopular* in DL?

First-order optimization

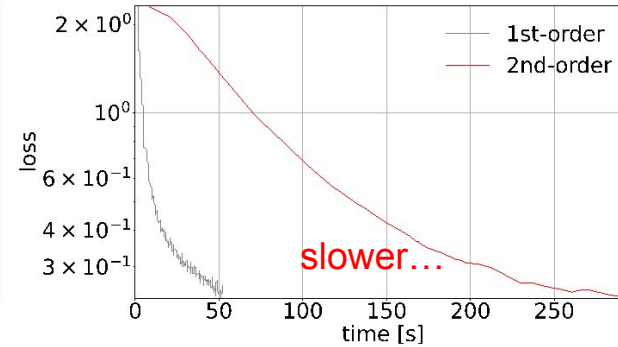
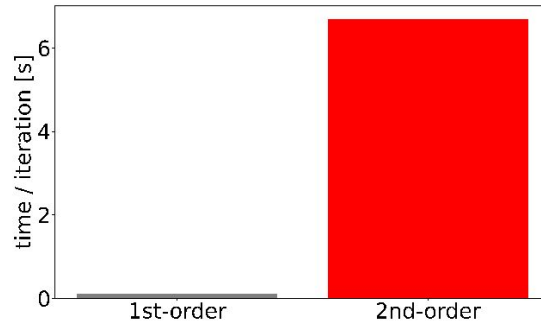
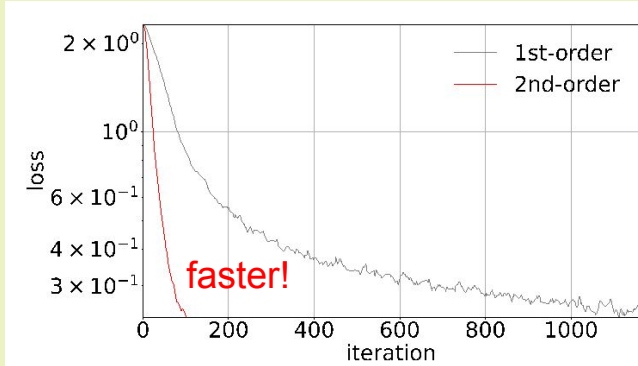
Second-order optimization

Forward/backward

$$\nabla \mathcal{L} \in \mathbb{R}^P$$

$$\nabla \mathcal{L} \in \mathbb{R}^P$$

Training a 3-layered MLP on MNIST



Overhead
 $\mathcal{O}(P^3)$

$\mathcal{L}(\theta^{(t)})$

Iterations to converge

×

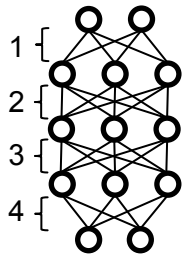
Time / iteration

=

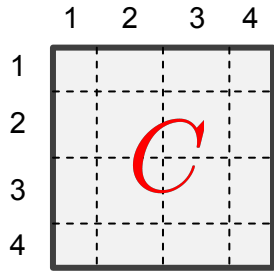
Total time

Kronecker-factored Approximate Curvature (K-FAC) Martens and Grosse, 2015

Neural network
(4 layers)

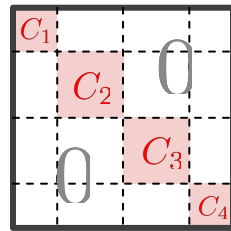


Full
(4 x 4 blocks)



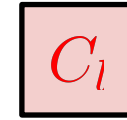
Step1. Layer-wise block-diagonal

Layer-wise
(4 blocks)



$$\mathcal{O}(P^3) \rightarrow \mathcal{O}(P_l^3)$$

Step2. Kronecker-factorization (for each layer)



Kronecker product

$$A \in \mathbb{R}^{d_1 \times d_2} \quad B \in \mathbb{R}^{d_3 \times d_4}$$

$$A \otimes B := \begin{bmatrix} A_{1,1}B & \dots & A_{1,d_2}B \\ \vdots & \ddots & \vdots \\ A_{d_1,1}B & \dots & A_{d_1,d_2}B \end{bmatrix} \in \mathbb{R}^{d_1 d_3 \times d_2 d_4}$$

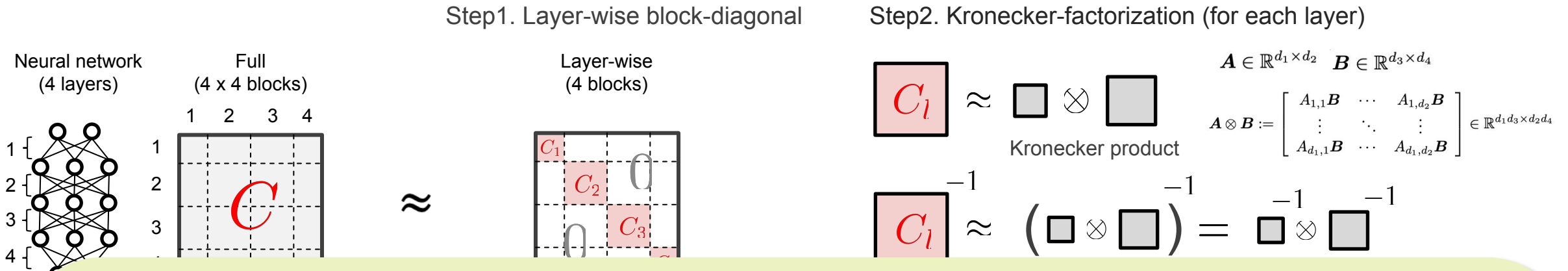


$$\left(\square \otimes \square \right)^{-1} = \square^{-1} \otimes \square^{-1}$$

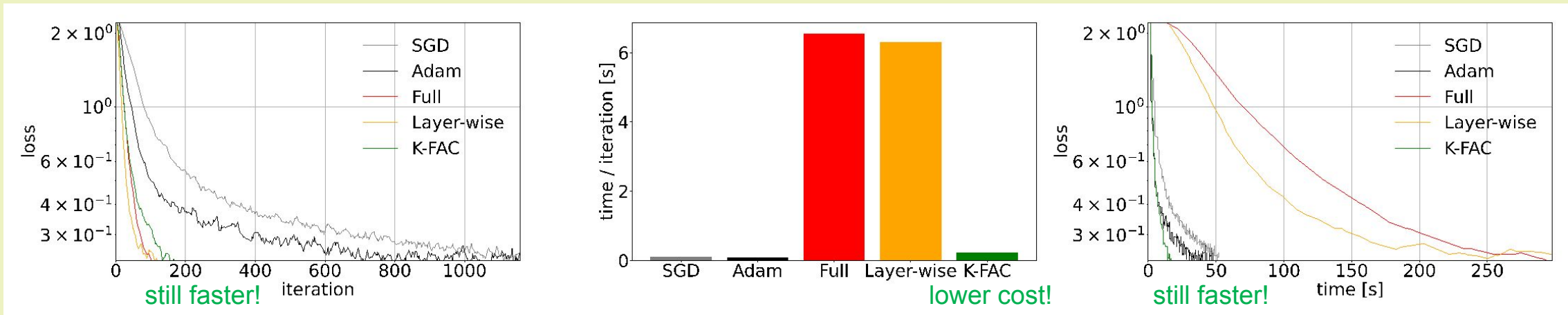
$$\mathcal{O}(P_l^3) \rightarrow \mathcal{O}(P_l^{3/2})$$

Kronecker-factored Approximate Curvature (K-FAC)

Martens and Grosse, 2015



Training a 3-layered MLP ($P_1 = 12560, P_2 = 136, P_3 = 90$) on MNIST



Iterations to converge

×

Time / iteration

=

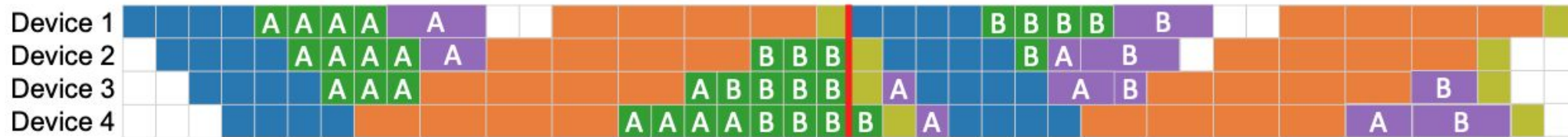
Total time

PipeFisher

(a) GPipe



(b) PipeFisher for GPipe



Pipeline bubble — Pipeline flush @ Device 1

What exactly are the work of K-FAC?

Work of K-FAC

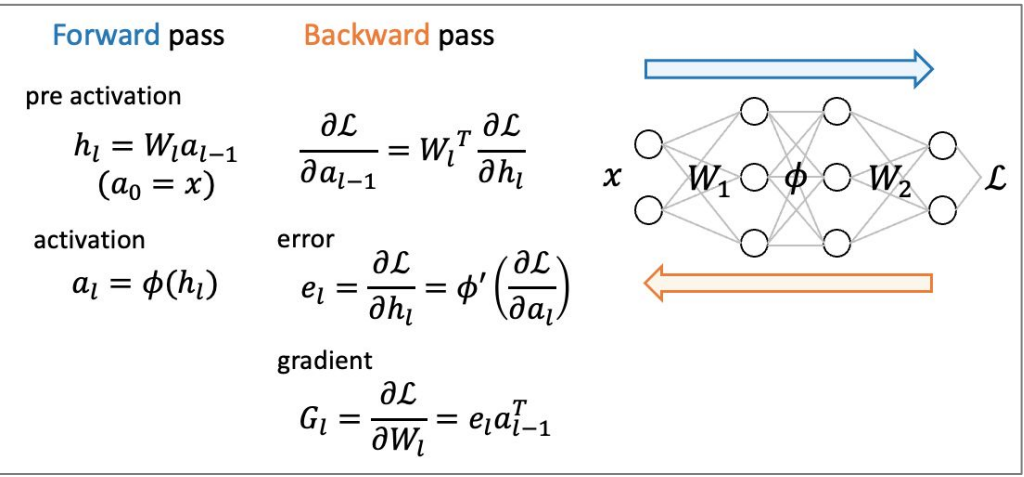
(mini-batch) SGD

$$G_l = \frac{\partial \mathcal{L}}{\partial W_l} = U_{e,l} \cdot U_{a,l-1}^\top$$

$(d_l^{out} \times B) (B \times d_l^{in})$ B : mini-batch size
 ← batched error ← batched activation

$$g_l = \text{vec}(G_l) \in \mathbb{R}^{P_l}$$

$$P_l = d_l^{out} \cdot d_l^{in}$$



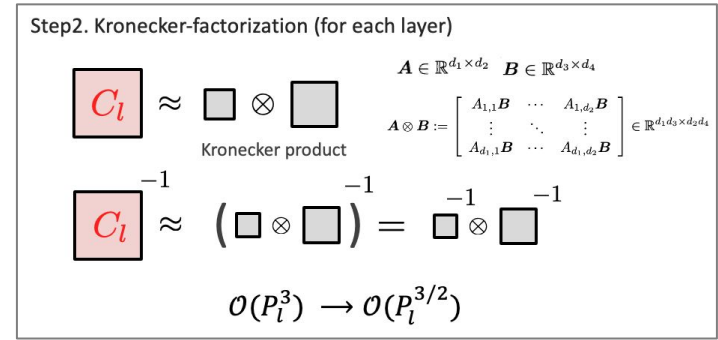
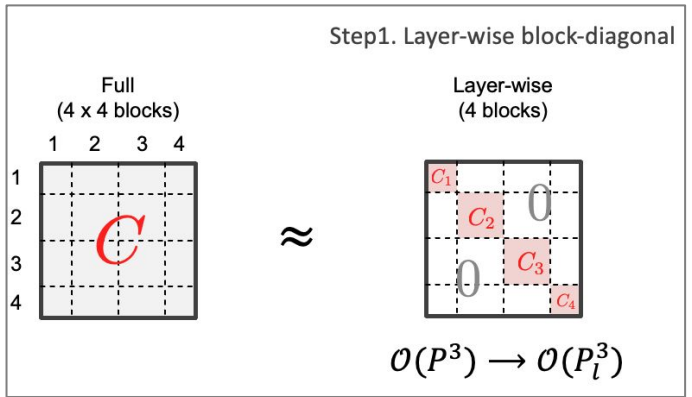
K-FAC

$$g_l^{kfac} \approx F_l^{-1} g_l$$

Fisher block for l-th layer

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta C^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

Precondition the gradient by the **curvature matrix**



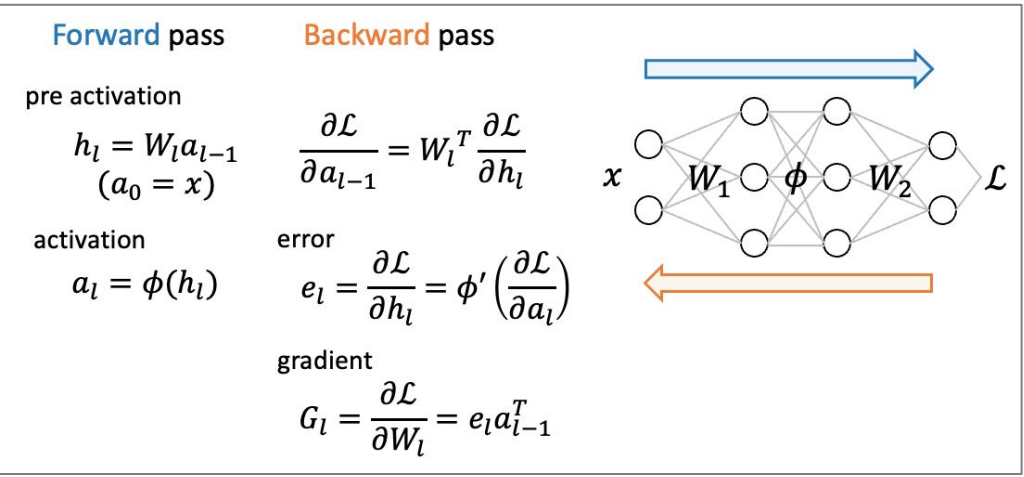
Work of K-FAC

(mini-batch)
SGD

$$G_l = \frac{\partial \mathcal{L}}{\partial W_l} = U_{e,l} \cdot U_{a,l-1}^\top$$

$(d_l^{out} \times B) \quad (B \times d_l^{in}) \quad B : \text{mini-batch size}$

$g_l = \text{vec}(G_l) \in \mathbb{R}^{P_l}$
 $P_l = d_l^{out} \cdot d_l^{in}$



K-FAC $g_l^{kfac} \approx F_l^{-1} g_l \approx (A_l \otimes B_l)^{-1} g_l = (A_l^{-1} \otimes B_l^{-1}) g_l = \text{vec}(B_l^{-1} G_l A_l^{-1})$

Curvature work

$A = U_a @ U_a.T$
 $B = U_e @ U_e.T$

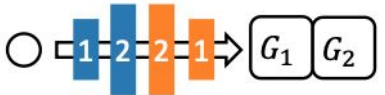
Inversion work

$A_{inv} = \text{inverse}(A)$
 $B_{inv} = \text{inverse}(B)$

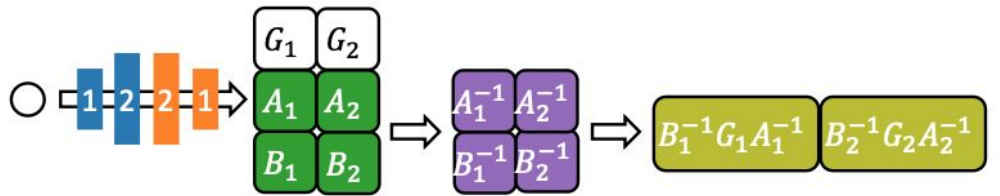
Precondition work

$G_{kfac} = B_{inv} @ G @ A_{inv}$

(a) SGD



(b) K-FAC

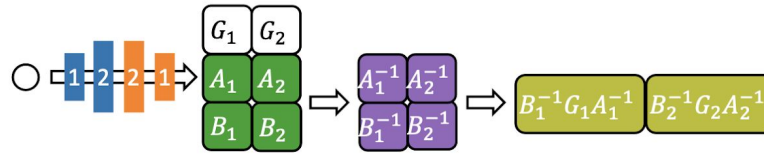
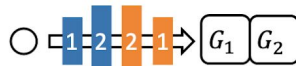


K-FAC with parallelism

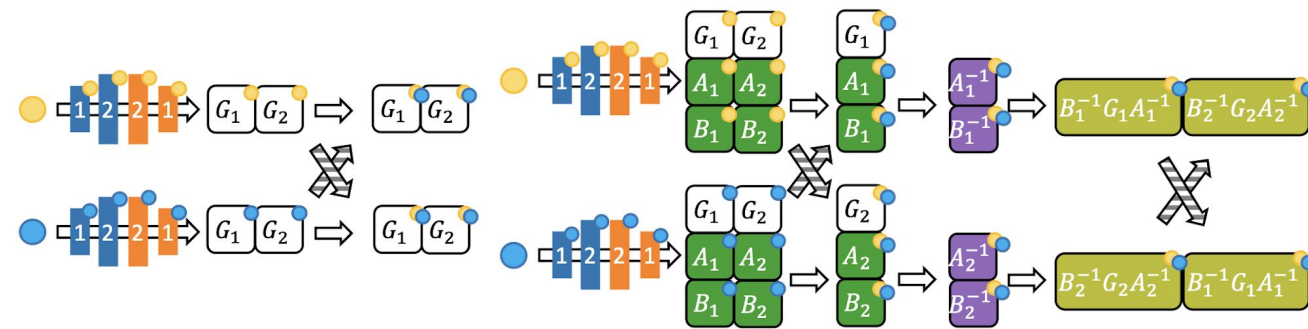
(a) SGD

(b) K-FAC

(i) No parallelism
(w/ 1 device)

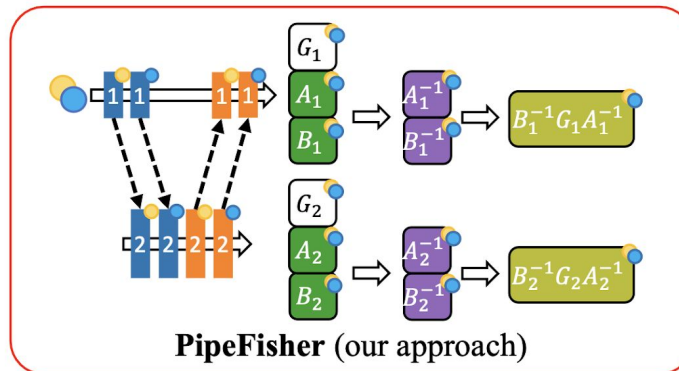
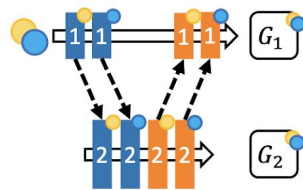





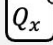


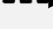
(ii) Data parallelism
(w/ 2 devices)



(+ inversion parallelism) Osawa et al., 2019

(iii) Pipeline parallelism
(w/ 2 devices)

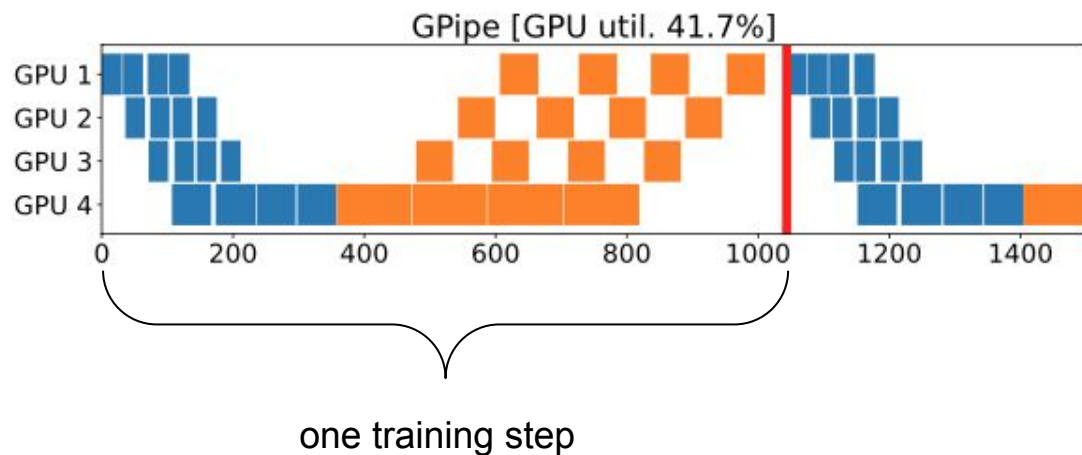


-  Mini- or micro-batch
-  Forward/backward at x-th layer for 
-  Quantity for x-th layer calculated for 
-  Collective communication
-  Point-to-point communication (send/recv)

1. Less memory consumption.
2. Inversion work are split without collective communication.
3. Better accelerator utilization.

Profiled results (1/2) [GPipe (Huang et al, 2018)]

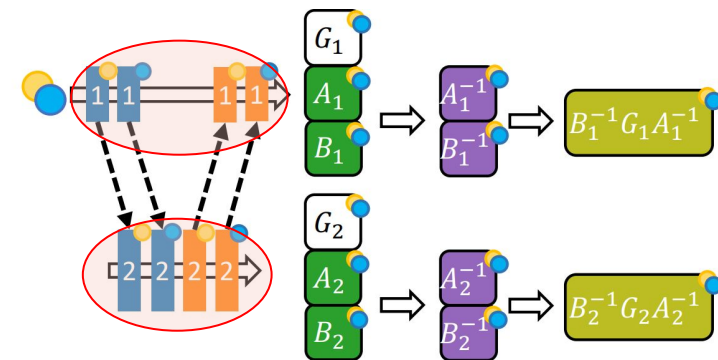
- BERT-Base (12 Transformer layers) w/ 4 pipeline stages (3 layers per stage) and 4 micro-batches
- CUDA kernel execution times on NVIDIA P100 GPUs



— flush @ GPU1 ■ forward ■ backward

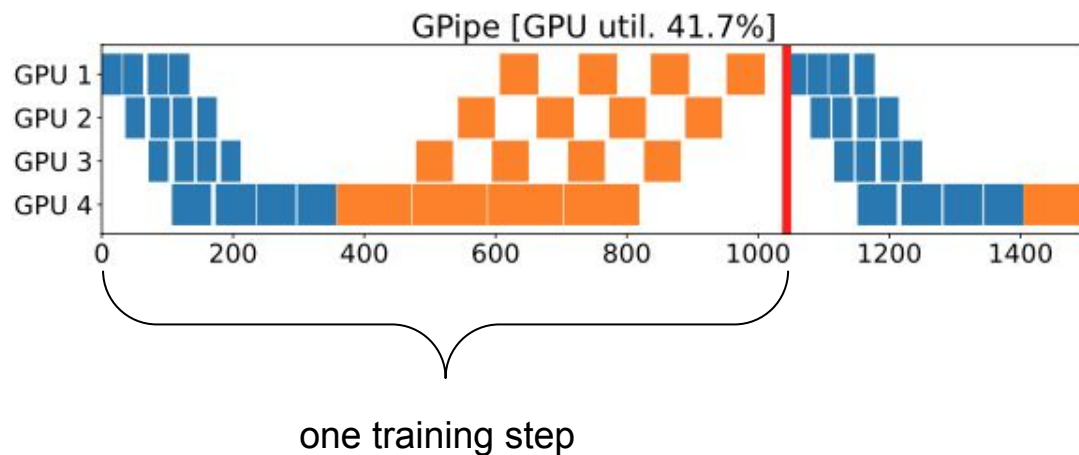
PipeFisher

- Step 1: **Measure** the times for the **forward/backward** works and **bubbles** in a training step.



Profiled results (1/2) [GPipe (Huang et al, 2018)]

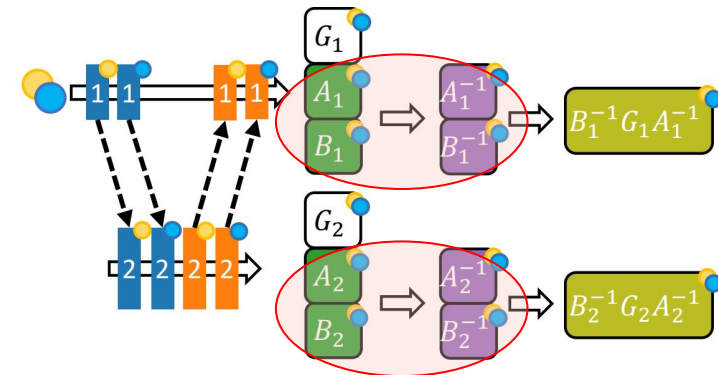
- BERT-Base (12 Transformer layers) w/ 4 pipeline stages (3 layers per stage) and 4 micro-batches
- CUDA kernel execution times on NVIDIA P100 GPUs



— flush @ GPU1 — forward — backward

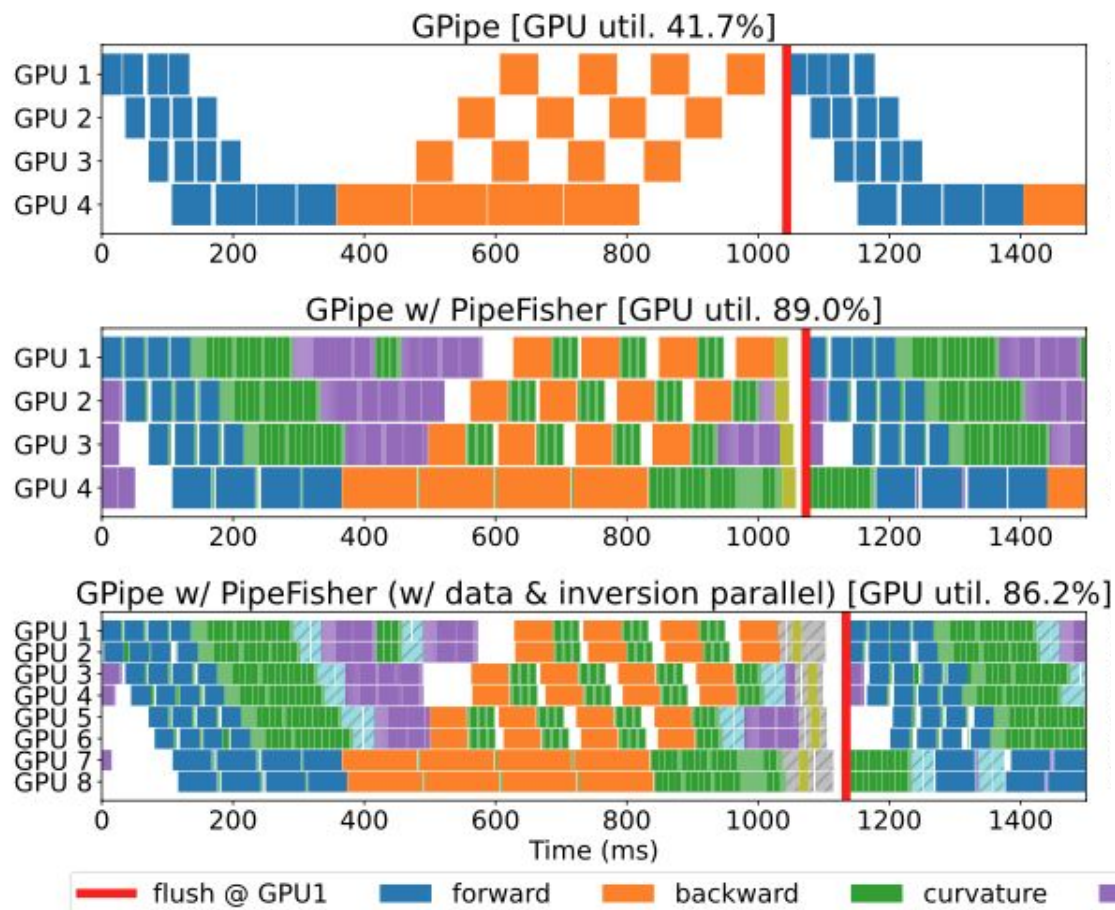
PipeFisher

- ❑ Step 1: **Measure** the times for the **forward/backward** works and **bubbles** in a training step.
- ❑ Step 2: **Measure** the times for the **curvature/inverse** works in a training step.



Profiled results (1/2) [GPipe (Huang et al, 2018)]

- BERT-Base (12 Transformer layers) w/ 4 pipeline stages (3 layers per stage) and 4 micro-batches
- CUDA kernel execution times on NVIDIA P100 GPUs



PipeFisher

- ❑ Step 1: **Measure** the times for the **forward/backward** works and **bubbles** in a training step.
- ❑ Step 2: **Measure** the times for the **curvature/inverse** works in a training step.
- ❑ Step 3: **Assign** the **curvature/inverse** works to the **bubbles** within a training step(s) and the **precondition** work at the end of every step.

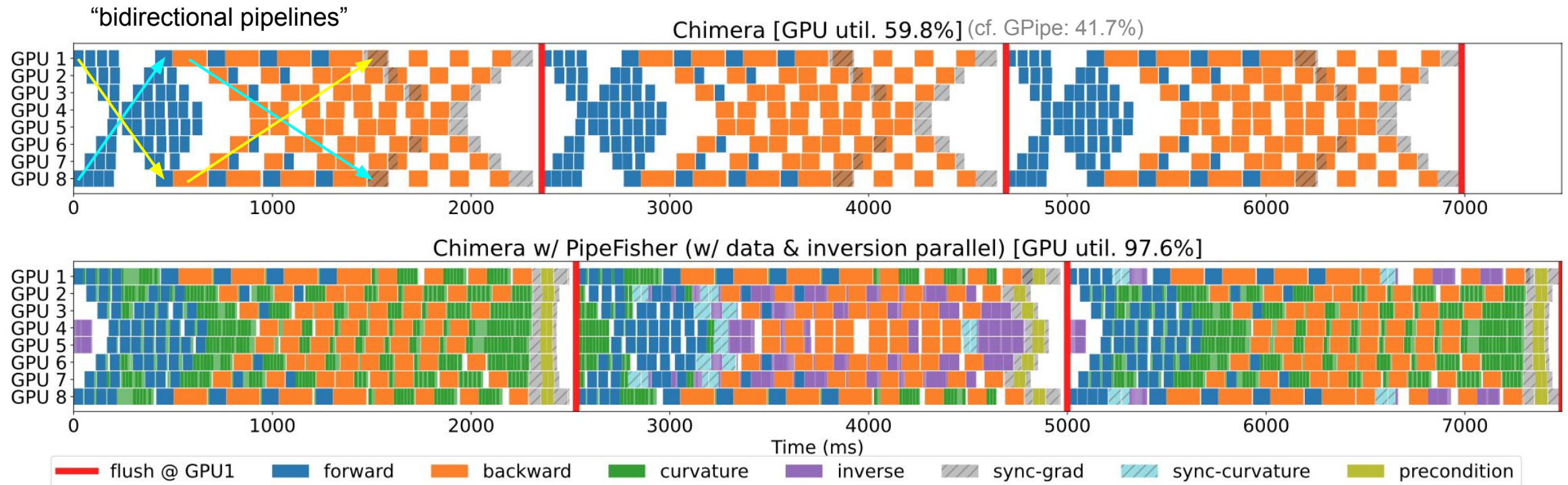
⚠ the ratio = (curvature+inverse) / bubbles determines the frequency of updating the preconditioner (e.g., every 1-2 steps).

$$g_l^{kfac} \approx F_l^{-1} g_l$$

10-100x higher freq. than common practice!

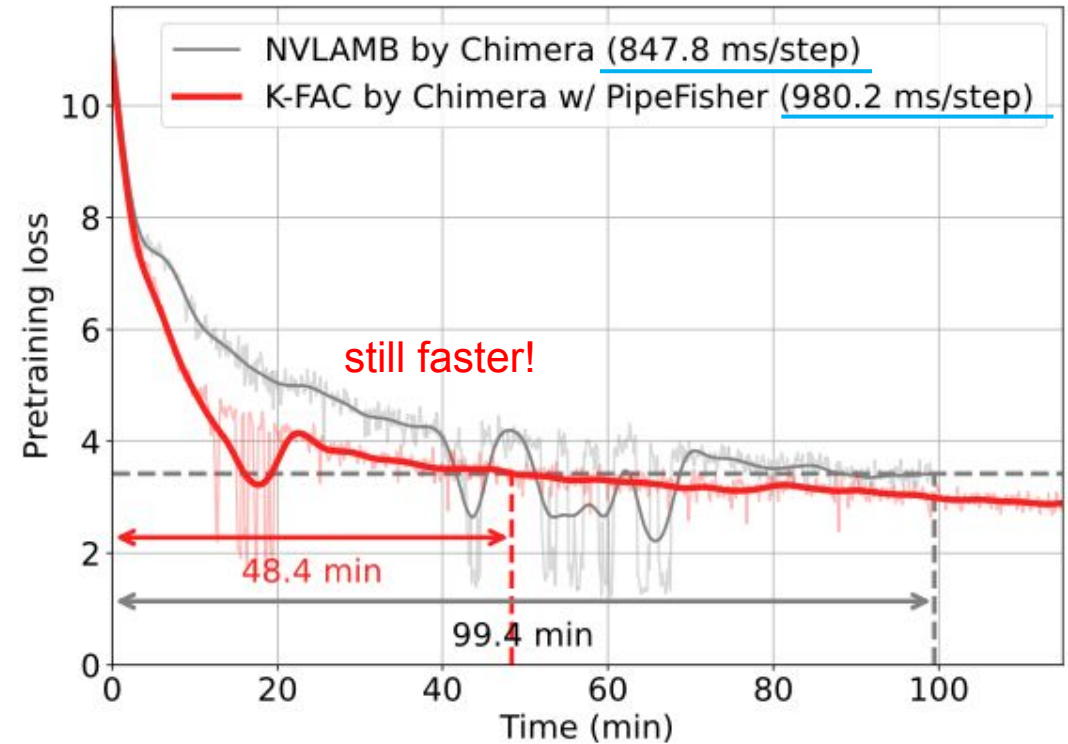
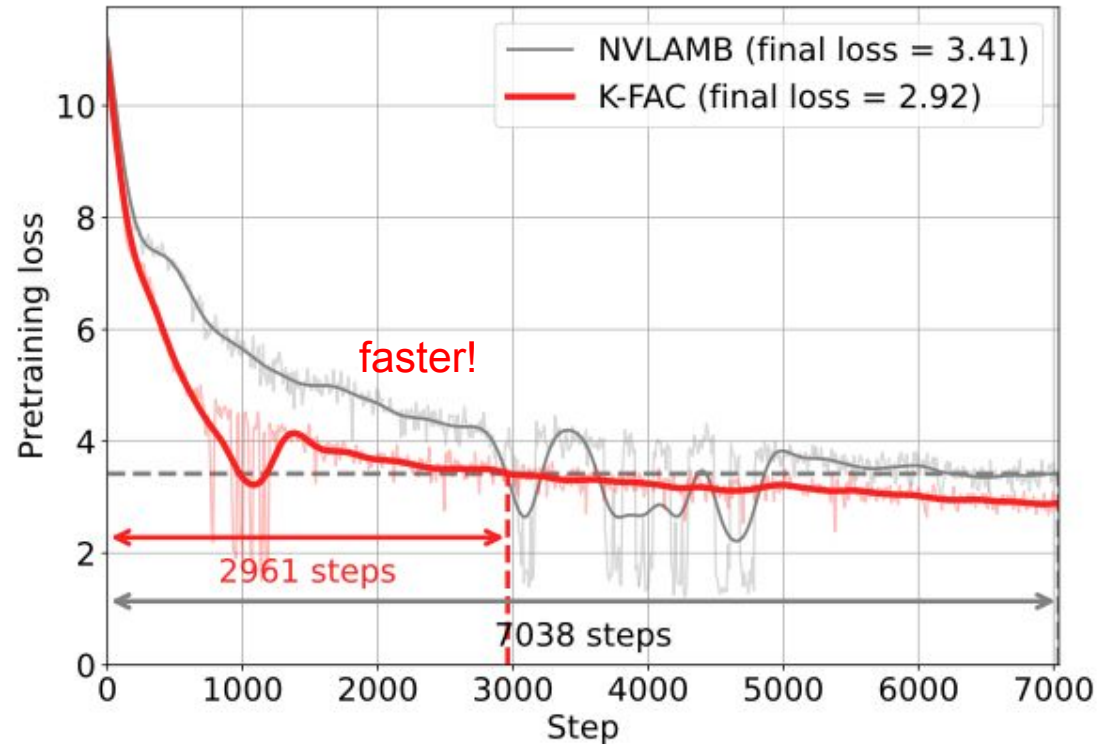
Profiled results (2/2) [Chimera (Li and Hoefler, 2021)]

- BERT-Large (24 Transformer layers) w/ 8 pipeline stages (3 layers per stage) and 8 micro-batches
- CUDA kernel execution times on NVIDIA P100 GPUs



BERT-Base Pretraining

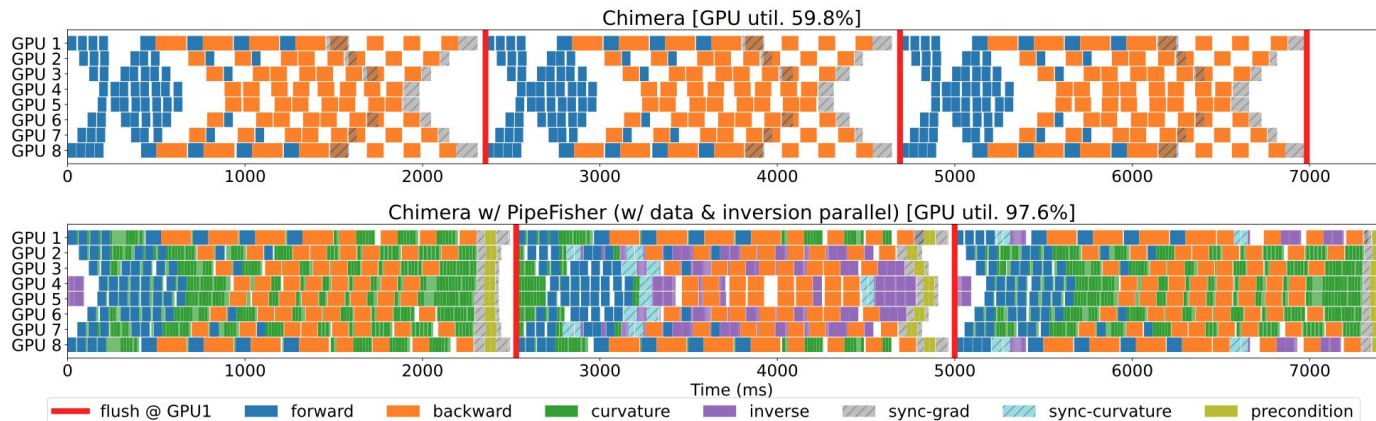
- BERT-Base (12 Transformer layers) w/ 4 pipeline stages (3 layers per stage) and 4 micro-batches
- Pretraining on the English Wikipedia
- Time measured on 256 NVIDIA P100 GPUs



BERT-Large Pretraining

- BERT-Large (24 Transformer layers) w/ 8 pipeline stages (3 layers per stage) and 8 micro-batches
- Pretraining on the English Wikipedia
- Time measured on 8 NVIDIA P100 GPUs (total training time is simulated)

Optimizer	Pipeline scheme	Phase 1			Phase 2	F1
		Steps	Time/step*	Time*	Steps	
NVLAMB	Chimera	7038	2345.6 ms	275.1 min	1563	90.1%
K-FAC	Chimera w/ PipeFisher	5000	2499.5 ms	208.3 min	1563	90.15%

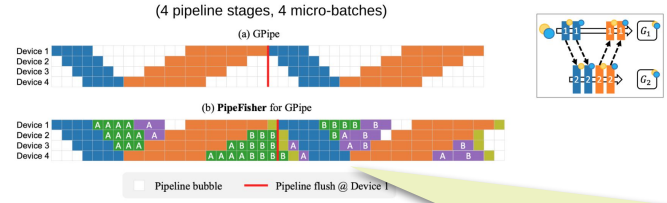



Conclusions

Bubbles in pipeline

- Pipelining creates **bubbles** of time in which accelerators become idle.
- The overhead of pipelining mainly comes from the **low utilization of accelerators**
- We suggest to assign **extra work** to the bubbles to gain **auxiliary benefits**.
- **PipeFisher** automatically assigns the **work of K-FAC** (a second-order optimization method based on the Fisher information matrix) to the bubbles for **accelerating training**.

(4 pipeline stages, 4 micro-batches)



 We can fill pipeline bubbles with **extra work**!

Conclusions

Bubbles in pipeline

- Pipelining creates **bubbles** of time in which accelerators become idle.
- The overhead of pipelining mainly comes from the **low utilization of accelerators**
- We suggest to assign **extra work** to the bubbles to gain **auxiliary benefits**.
- PipeFisher** automatically assigns the **work of K-FAC** (a second-order optimization method based on the Fisher information matrix) to the bubbles for **accelerating training**.

(4 pipeline stages, 4 micro-batches)

(a) GPipe

(b) PipeFisher for GPipe

Legend: Pipeline bubble (grey square), Pipeline flush @ Device 1 (red line)

Why second-order optimization?

First-order Optimization (gradient descent)


$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)})$$

Second-order Optimization

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \mathbf{C}^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

Precondition the gradient by the **curvature matrix**

Figure from J. Martens, 2010

 **Faster convergence (loss vs # steps) !**

Conclusions

Bubbles in pipeline

- Pipelining creates **bubbles** of time in which accelerators become idle.
- The overhead of pipelining mainly comes from the **low utilization of accelerators**
- We suggest to assign **extra work** to the bubbles to gain **auxiliary benefits**.
- PipeFisher** automatically assigns the **work of K-FAC** (a second-order optimization method based on the Fisher information matrix) to the bubbles for **accelerating training**.

(4 pipeline stages, 4 micro-batches)

Why second-order optimization?

First-order Optimization (gradient descent)

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)})$$

Second-order Optimization

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \mathbf{C}^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

Precondition the gradient by the **curvature matrix**

Figure from J. Martens, 2010

K-FAC with parallelism

(i) No parallelism (w/ 1 device)

(ii) Data parallelism (w/ 2 devices)

(iii) Pipeline parallelism (w/ 2 devices)

PipeFisher (our approach)

Legend:

- Mini- or micro-batch
- ↔ Forward/backward at x-th layer for ○
- Quantity for x-th layer calculated for ○
- ⊗ Collective communication
- Point-to-point communication (end-to-end)

💡

- **PipeFisher** -> K-FAC in pipeline bubbles
- K-FAC (layer-wise preconditioning) is compatible with pipelining!

Conclusions

Bubbles in pipeline

- Pipelining creates **bubbles** of time in which accelerators become idle.
- The overhead of pipelining mainly comes from the **low utilization of accelerators**
- We suggest to assign **extra work** to the bubbles to gain **auxiliary benefits**.
- PipeFisher** automatically assigns the **work of K-FAC** (a second-order optimization method based on the Fisher information matrix) to the bubbles for **accelerating training**.

K-FAC with parallelism

Why second-order optimization?

First-order Optimization (gradient descent)

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)})$$

Second-order Optimization

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \mathbf{C}^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

Precondition the gradient by the **curvature matrix**

Figure from J. Martens, 2010

- Higher GPU utilization!
- Faster training (F1 vs time) (reduction to 50-70%)!

BERT-Large Pretraining

- BERT-Large (24 Transformer layers) w/ 8 pipeline stages (3 layers per stage) and 8 micro-batches
- Pretraining on the English Wikipedia
- Time measured on 8 NVIDIA P100 GPUs (total training time is simulated)

Optimizer	Pipeline scheme	Phase 1			Phase 2 Steps	F1
		Steps	Time/step*	Time*		
NVLAMB	Chimera	7038	2345.6 ms	275.1 min	1563	90.1%
K-FAC	Chimera w/ PipeFisher	5000	2499.5 ms	208.3 min	1563	90.15%

Conclusions

Bubbles in pipeline

- Pipelining creates **bubbles** of time in which accelerators become idle.
- The overhead of pipelining mainly comes from the **low utilization of accelerators**
- We suggest to assign **extra work** to the bubbles to gain **auxiliary benefits**.
- PipeFisher** automatically assigns the **work of K-FAC** (a second-order optimization method based on the Fisher information matrix) to the bubbles for **accelerating training**.

(4 pipeline stages, 4 micro-batches)

(a) GPIPE

(b) PipeFisher for GPIPE

Legend: Pipeline bubble (grey), Pipeline flush @ Device 1 (red)

K-FAC with parallelism

(i) No parallelism (w/ 1 device)

(ii) Data parallelism (w/ 2 devices)

(iii) Pipeline parallelism (w/ 2 devices)

Legend: Mini- or micro-batch (circle), Forward/backward at x-th layer for (square), Quantity for x-th layer calculated for (circle), Collective communication (cross), Point-to-point communication (end/recv) (dashed line)

Chimera w/ PipeFisher (w/ data)

Legend: Sub-gp GPIPE (red), forward (blue), backward (orange), compute (green)

Why second-order optimization?

First-order Optimization (gradient descent)

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)})$$

Second-order Optimization

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \mathbf{C}^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

Precondition the gradient by the **curvature**

BERT-Large Pretraining

- BERT-Large (24 Transformer layers) w/ 8 pipeline
- Pretraining on the English Wikipedia
- Time measured on 8 NVIDIA P100 GPUs (total)

Optimizer	Pipeline scheme	Steps
NVLAMB	Chimera	7038
K-FAC	Chimera w/ PipeFisher	5000

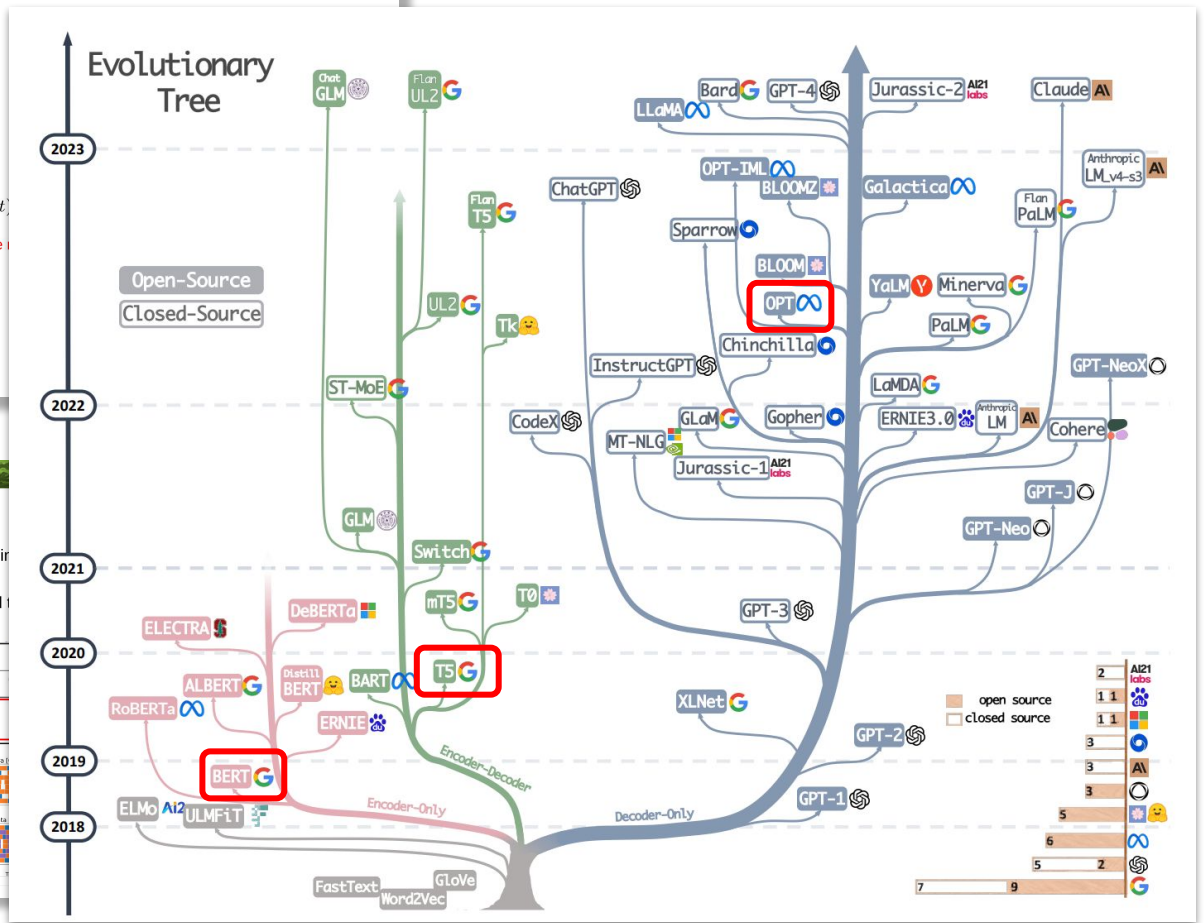


Image: Yang et al., "Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond", 2023.

Conclusions

Bubbles in pipeline

- Pipelining creates **bubbles** of time in which accelerators become idle.
- The overhead of pipelining mainly comes from the **low utilization of accelerators**
- We suggest to assign **extra work** to the bubbles to gain **auxiliary benefits**.
- PipeFisher** automatically assigns the **work of K-FAC** (a second-order optimization method based on the Fisher information matrix) to the bubbles for **accelerating training**.

Why second-order optimization?

First-order Optimization (gradient descent)

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)})$$

Second-order Optimization

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \mathbf{C}^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

Precondition the gradient by the **curvature matrix**

K-FAC with parallelism

- (i) No parallelism (w/ 1 device)
- (ii) Data parallelism (w/ 2 devices)
- (iii) Pipeline parallelism (w/ 2 devices)
- PipeFisher (our approach)**

BERT-Large Pretraining

- BERT-Large (24 Transformer layers) w/ 8 pipeline stages (3 layers per stage) and 8 micro-batches
- Pretraining on the English Wikipedia
- Time measured on 8 NVIDIA P100 GPUs (total training time is simulated)

Optimizer	Pipeline scheme	Phase 1			Phase 2 Steps	F1
		Steps	Time/step*	Time*		
NVLAMB	Chimera	7038	2345.6 ms	275.1 min	1563	90.1%
K-FAC	Chimera w/ PipeFisher	5000	2499.5 ms	208.3 min	1563	90.15%

More of SPCL's research:

- [youtube.com/@spcl](https://www.youtube.com/@spcl) **180+ Talks**
- twitter.com/spcl_eth **1.3K+ Followers**
- github.com/spcl **2K+ Stars**

... or spcl.ethz.ch



Paper link:

<https://arxiv.org/abs/2211.14133>

