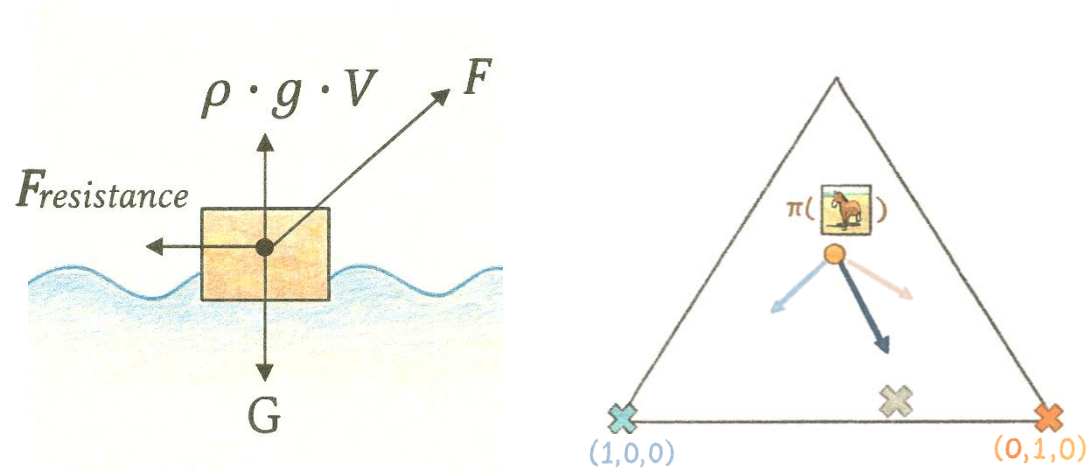# Learning Dynamics of Deep Learning
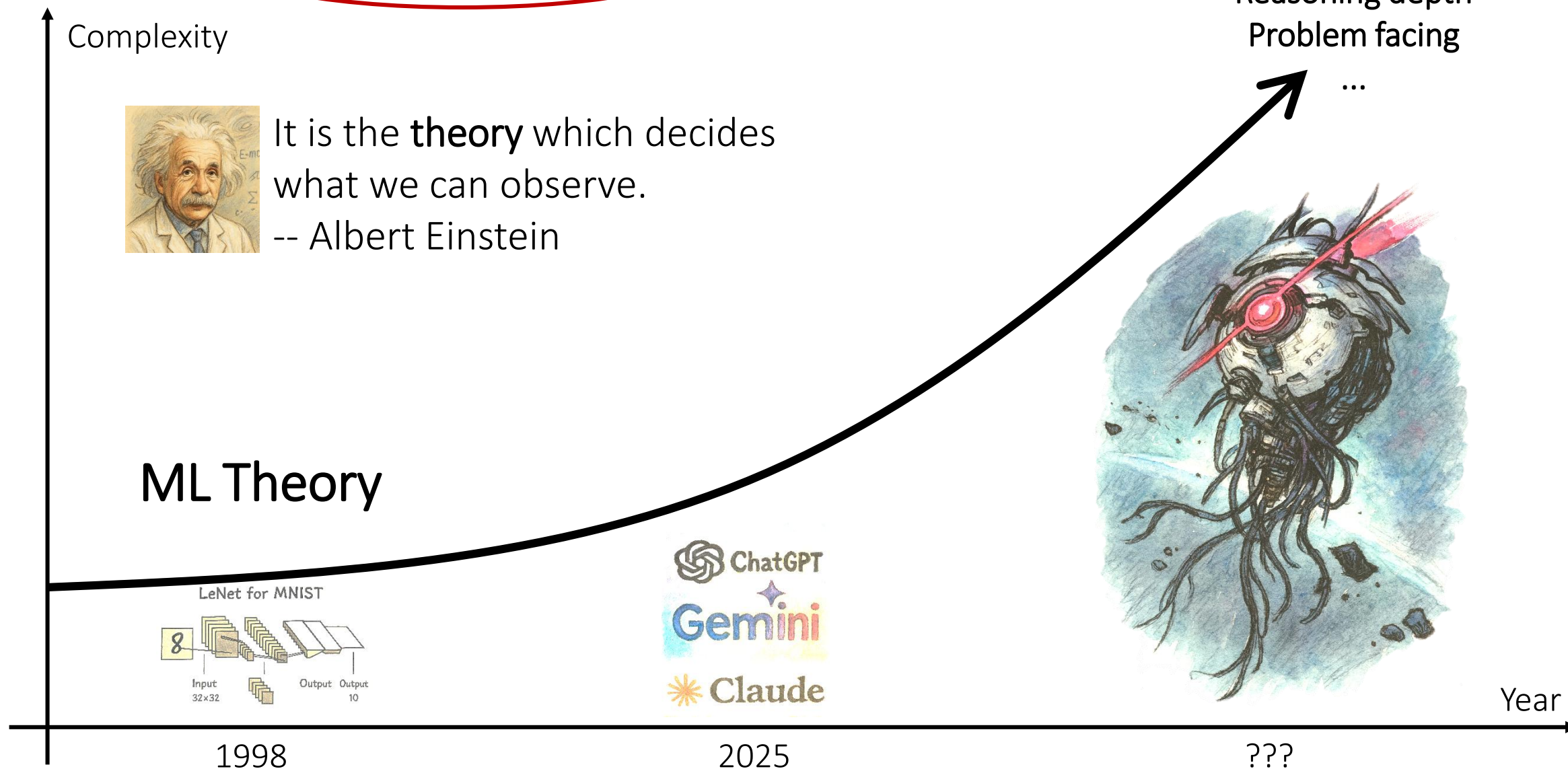
## -- Force Analysis of Deep Neural Networks



Yi (Joshua) Ren
Supervisor: Danica J. Sutherland
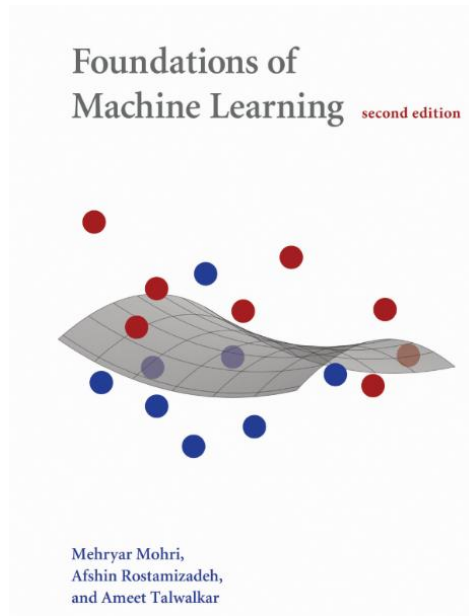
**Motivation: understanding and controlling AI systems need theory**



Model size
Reasoning depth
Problem facing
...

Complexity

It is the **theory** which decides
what we can observe.
-- Albert Einstein

ML Theory

LeNet for MNIST

Input
32×32

Output  Output
        10

ChatGPT

Gemini

Claude

Year

1998                    2025                    ???

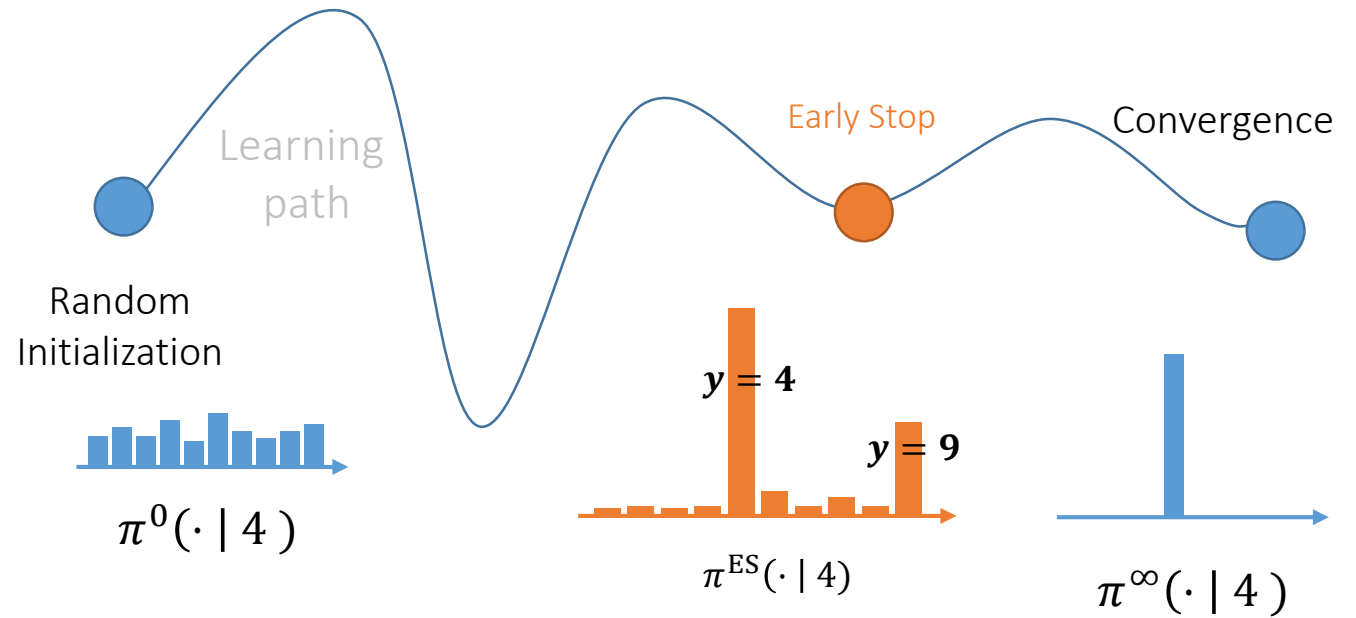# Motivation: ML theory needs diverse perspectives

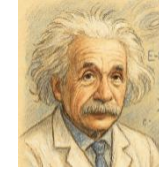- **PAC learning framework**
  -- strict, elegant, **global, and macroscopic**

- But, I failed to use it understanding this **emergent** behavior:
  -- an interesting pairing effect emerges during training

Foundations of
Machine Learning *second edition*

Mehryar Mohri,
Afshin Rostamizadeh,
and Ameet Talwalkar

Learning path

Random Initialization

$\pi^0(\cdot \mid 4)$

Early Stop

Convergence

$y = 4$

$y = 9$

$\pi^{ES}(\cdot \mid 4)$

$\pi^\infty(\cdot \mid 4)$
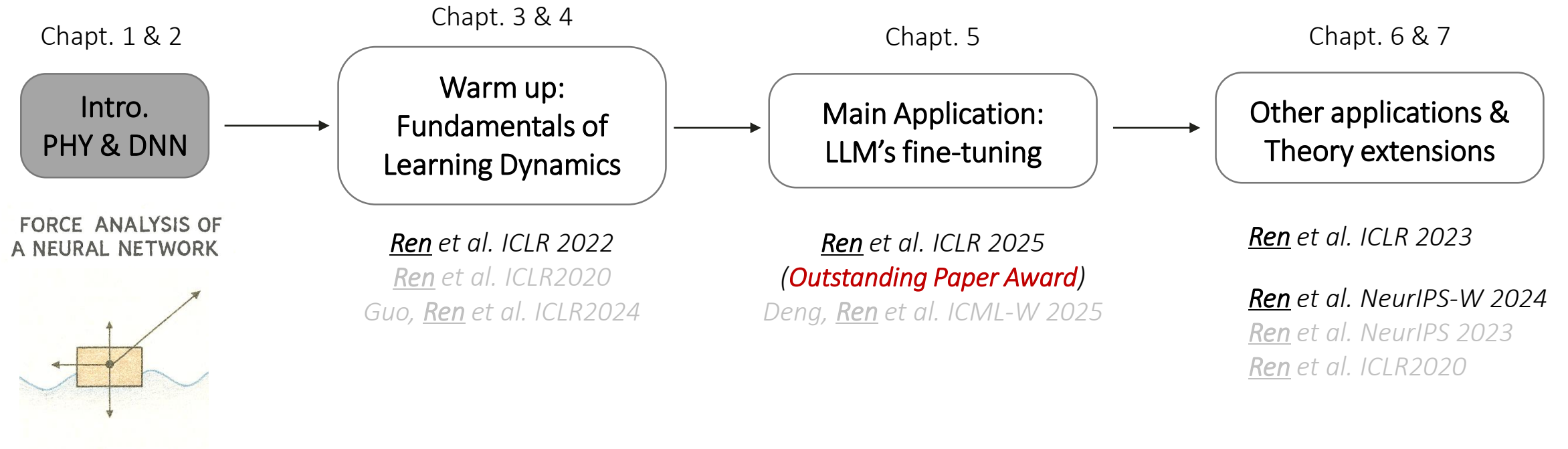
**Methodology: zoom in, in time and sample spaces**

It is the theory which decides what we can observe.
-- Albert Einstein

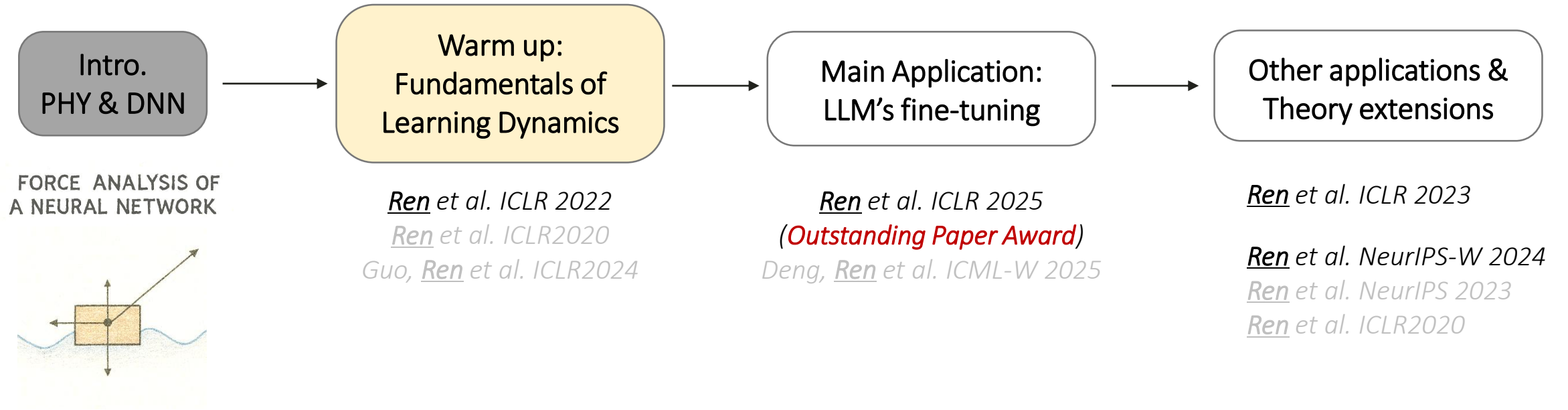# Force analysis of Neural Network (Learning Dynamics of Deep Learning)

## A fine-grained, physics-inspired ML theoretical framework

# Outline

Chapt. 1 & 2

Chapt. 3 & 4

Chapt. 5

Chapt. 6 & 7

**Intro.
PHY & DNN**

**Warm up:
Fundamentals of
Learning Dynamics**

**Main Application:
LLM's fine-tuning**

**Other applications &
Theory extensions**

FORCE ANALYSIS OF
A NEURAL NETWORK

*__Ren__ et al. ICLR 2022*
*__Ren__ et al. ICLR2020*
*Guo, __Ren__ et al. ICLR2024*

*__Ren__ et al. ICLR 2025*
*(Outstanding Paper Award)*
*Deng, __Ren__ et al. ICML-W 2025*

*__Ren__ et al. ICLR 2023*

*__Ren__ et al. NeurIPS-W 2024*
*__Ren__ et al. NeurIPS 2023*
*__Ren__ et al. ICLR2020*

# Outline



Intro.
PHY & DNN

FORCE ANALYSIS OF
A NEURAL NETWORK

Warm up:
Fundamentals of
Learning Dynamics

*Ren* et al. ICLR 2022
*Ren* et al. ICLR2020
Guo, *Ren* et al. ICLR2024

Main Application:
LLM's fine-tuning

*Ren* et al. ICLR 2025
*(Outstanding Paper Award)*
*Deng, Ren* et al. ICML-W 2025

Other applications &
Theory extensions

*Ren* et al. ICLR 2023

*Ren* et al. NeurIPS-W 2024
*Ren* et al. NeurIPS 2023
*Ren* et al. ICLR2020

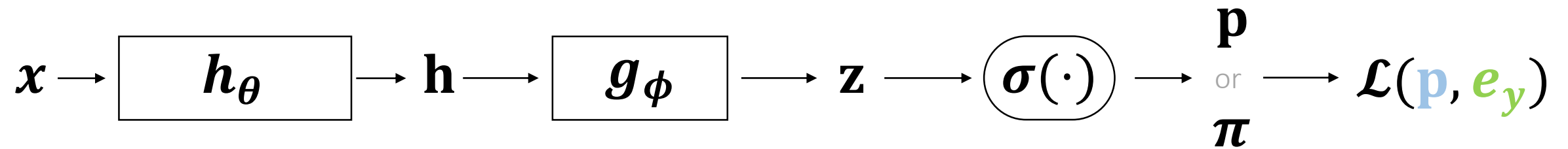# BETTER SUPERVISORY SIGNALS BY OBSERVING LEARNING PATHS

**Yi Ren**
UBC
renyi.joshua@gmail.com

**Shangmin Guo**
University of Edinburgh
s.guo@ed.ac.uk

**Danica J. Sutherland**
UBC and Amii
dsuth@cs.ubc.ca
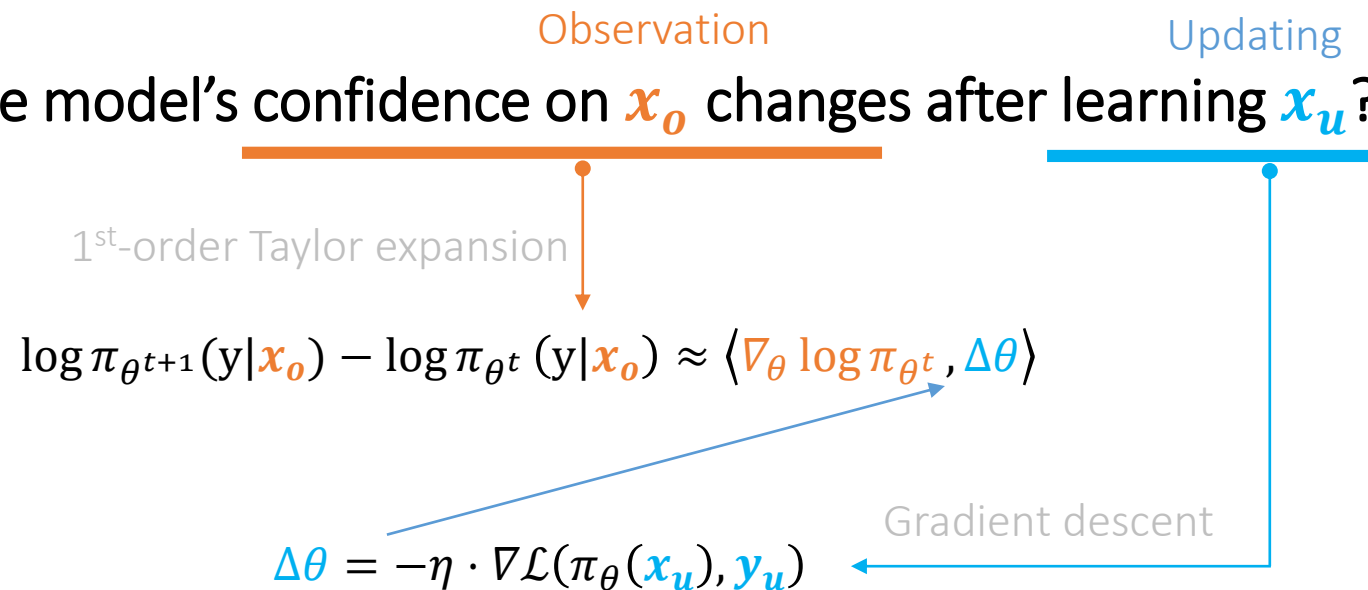
# Typical ML system: a sketch for notations

$$\mathcal{L}_{\text{ce}} = -\sum_{v=1}^{V} y_v \log\big(p(y = v|x)\big) = -\mathbf{e}_y^{\text{T}} \log \mathbf{p}(x) = -\mathbf{e}_y^{\text{T}} \log \boldsymbol{\sigma}(\mathbf{z}) = \cdots$$

$$x \rightarrow \boxed{\boldsymbol{h_\theta}} \rightarrow \mathbf{h} \rightarrow \boxed{\boldsymbol{g_\phi}} \rightarrow \mathbf{z} \rightarrow \big(\boldsymbol{\sigma(\cdot)}\big) \rightarrow \begin{matrix} \mathbf{p} \\ \text{or} \\ \boldsymbol{\pi} \end{matrix} \rightarrow \boldsymbol{\mathcal{L}(\mathbf{p}, e_y)}$$



$$\frac{e^{z_i}}{\sum_{j=1}^{V} e_j^z}$$

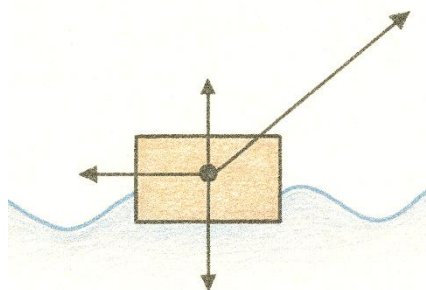| Input | Backbone | Hidden-embed  CLS Token | Task head  Can also be complex | Logits  Length-V vector | Softmax | Prediction  Length-V simplex | Loss  Depending on the task |

8

# Warm up: formalize the problem

Definition of one-step influence: **How the model's confidence on $x_o$ changes after learning $x_u$?**

- Analyze what?

  ✓ Model's prediction on $x_o$

- Where does the force comes from?

  ✓ Model's update on learning $x_u$

1st-order Taylor expansion

$$\log \pi_{\theta^{t+1}}(y|x_o) - \log \pi_{\theta^t}(y|x_o) \approx \langle \nabla_\theta \log \pi_{\theta^t}, \Delta\theta \rangle$$

Gradient descent

$$\Delta\theta = -\eta \cdot \nabla \mathcal{L}(\pi_\theta(x_u), y_u)$$

FORCE ANALYSIS OF A NEURAL NETWORK

$$\Delta \log \pi_{\theta^t}(y|x_o) = -\eta \mathcal{A}^t(x_o)\mathcal{K}^t(x_o, x_u)\mathcal{G}^t(x_u, y_u) + \mathcal{O}(\eta^2)$$

ICML 2017

PM LR Proceedings of Machine Learning Research
https://proceedings.mlr.press › ...   PDF

**Understanding Black-box Predictions via Influence Functions**

by PW Koh · Cited by 3508 — In this paper, we use **influence func- tions** — a classic technique from robust statis- tics — to trace a model's prediction through the learning algorithm and ...

9

**Warm up: understand the role of K-term**

$$\Delta \log \pi_{\theta^t}(y|\boldsymbol{x_o}) = -\eta \mathcal{A}^t(\boldsymbol{x_o}) \mathcal{K}^t(\boldsymbol{x_o}, \boldsymbol{x_u}) \mathcal{G}^t(\boldsymbol{x_u}, \boldsymbol{y_u}) + \mathcal{O}(\eta^2)$$

$$\nabla_{\mathbf{z}} \log \pi_{\theta^t} = I - \mathbf{1}(\pi^t)^\top = \begin{bmatrix} 1-\pi_1 & -\pi_1 & \cdots & -\pi_1 \\ -\pi_2 & 1-\pi_2 & \cdots & -\pi_2 \\ \cdots & \cdots & \ddots & \cdots \\ -\pi_V & -\pi_V & \cdots & 1-\pi_V \end{bmatrix}$$

Inner product of gradients
Empirical NTK
$$\nabla_\theta \boldsymbol{z_o}(\nabla_\theta \boldsymbol{z_u})^T$$
$$\pi = \text{Softmax}(\boldsymbol{z}); \; z = h_\theta(\boldsymbol{x})$$

$$\nabla_{\mathbf{z}} \mathcal{L}(\boldsymbol{x_u}, \boldsymbol{y_u})\Big|_{\mathbf{z}^t}$$
For cross-entropy
$$\pi_\theta(\boldsymbol{y}|\boldsymbol{x_u}) - \boldsymbol{e_{y_u}}$$

## Let's Warm up with a MNIST classification problem



Learn a "4" in this update

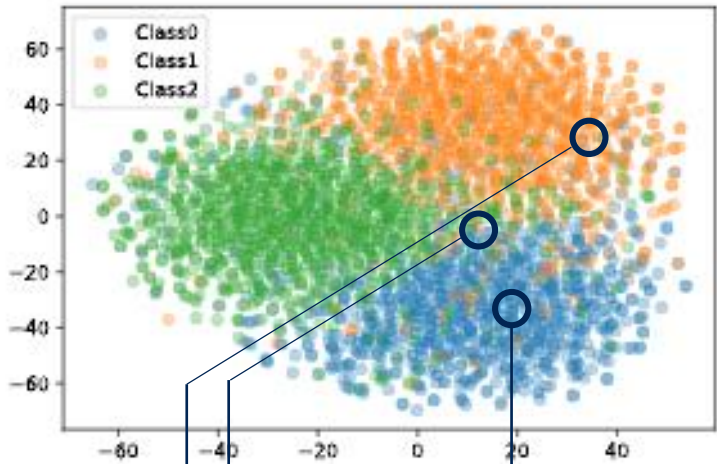Accumulates over several epochs    Imposed on $\boldsymbol{x_o}$    Projected by $\mathcal{K}^t$ Normalized by $\mathcal{A}^t$    Force from $\mathcal{G}^t$

# Warm up: understand the evolution of G-term

$$\Delta \log \pi_{\theta^t}(y|\textcolor{orange}{x_o}) \approx -\eta \mathcal{A}^t(\textcolor{orange}{x_o}) \, \mathcal{K}^t(\textcolor{orange}{x_o}, \textcolor{cyan}{x_u}) \, \mathcal{G}^t(\textcolor{cyan}{x_u}, \textcolor{cyan}{y_u})$$

- Examples with different difficulty



- Consider noisy-CIFAR-3
  (Numbers are sample ID)

All labeled as "Plane"

Easy:  $p^*(y|x) = [0.9, 0.1, 0]$
$\mathbf{e}_{y_n}^{\mathrm{T}} = [\mathbf{1}, 0, 0]$

A plane.

Medium:  $p^*(y|x) = [0.5, 0.3, 0.2]$
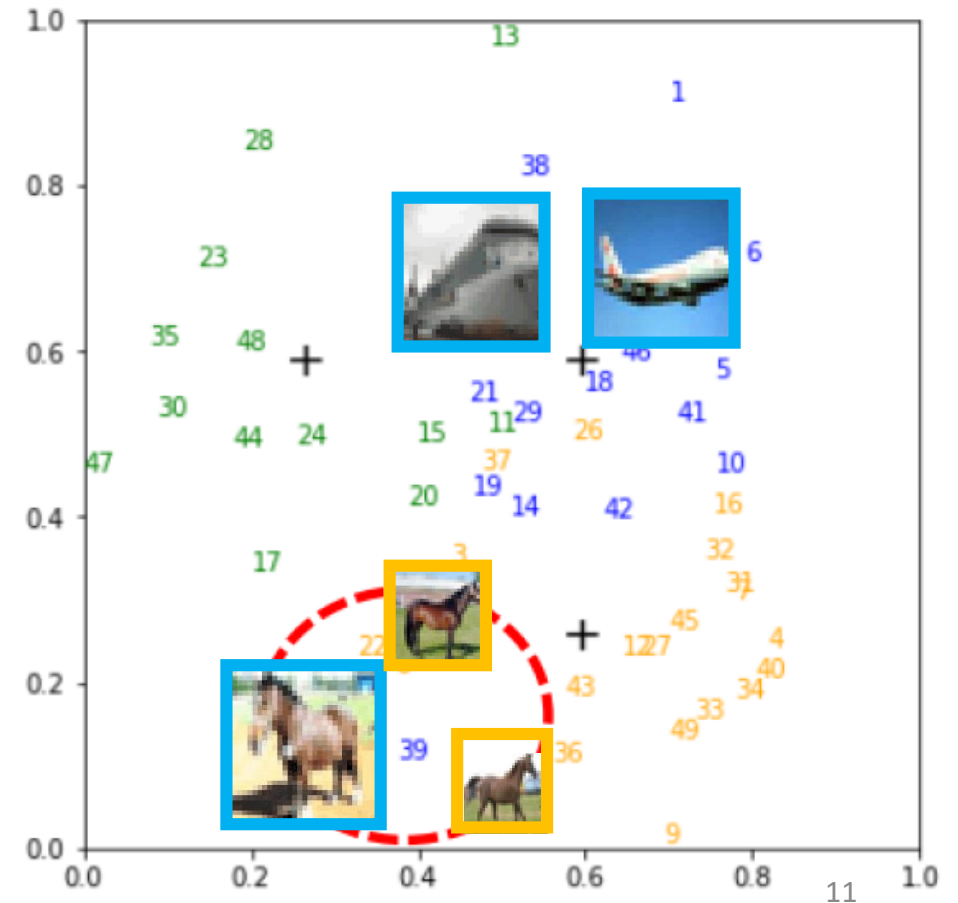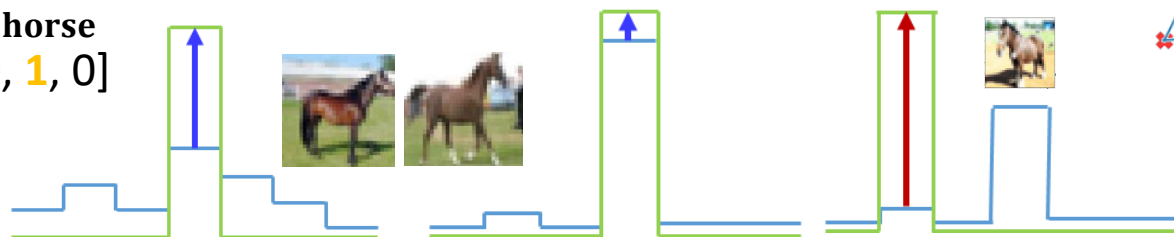$\mathbf{e}_{y_n}^{\mathrm{T}} = [\mathbf{1}, 0, 0]$

Plane? Ship?

Hard:
(Wrong label)  $p^*(y|x) = [0.1, 0.1, 0.8]$
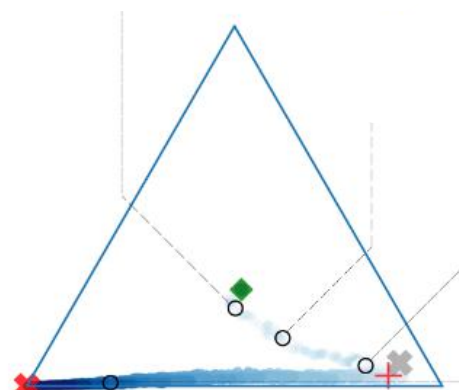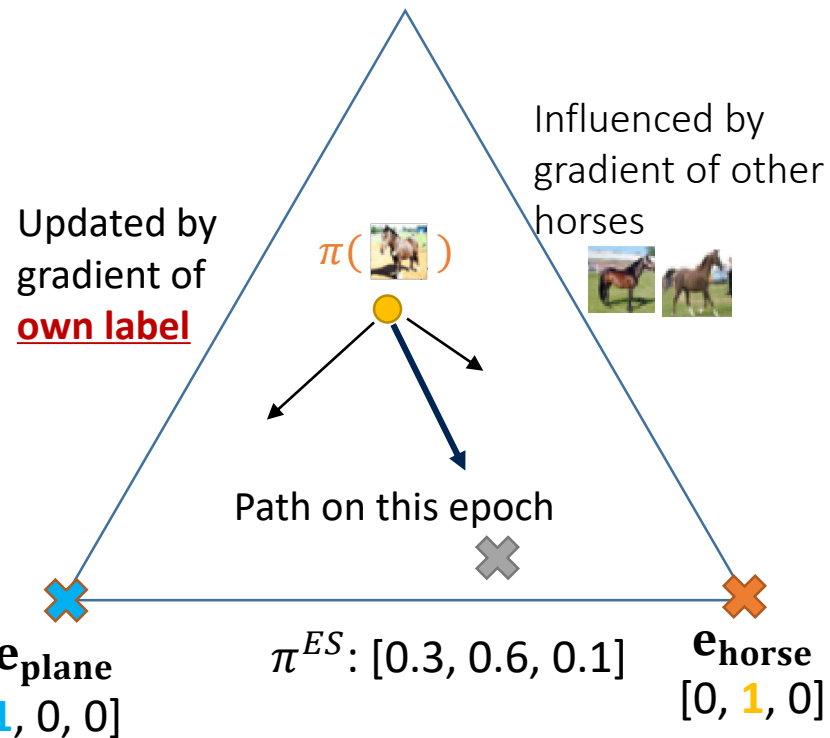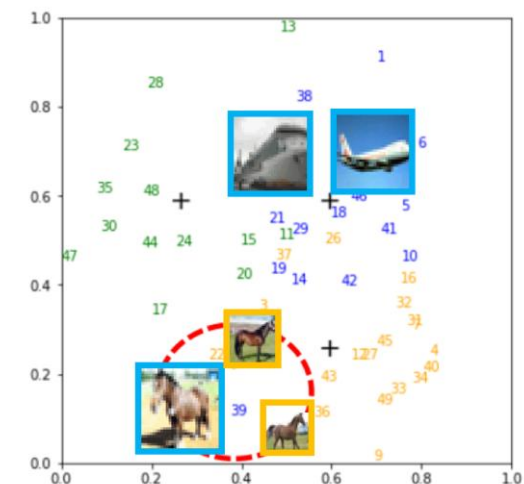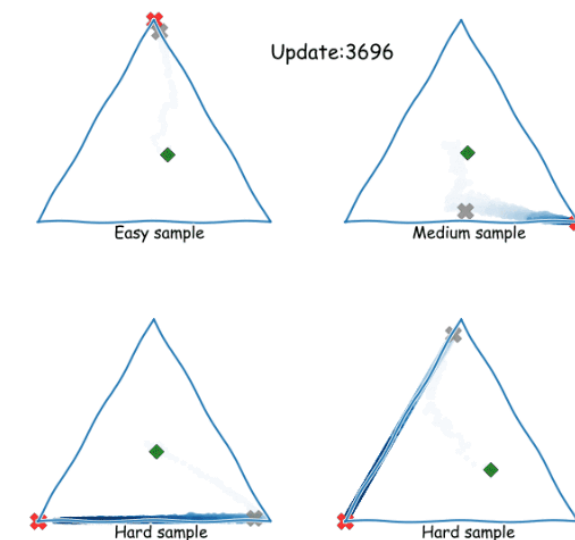$\mathbf{e}_{y_n}^{\mathrm{T}} = [\mathbf{1}, 0, 0]$

Plane????

# Warm up: understand the evolution of G-term

$$\Delta \log \pi_{\theta^t}(y|\boldsymbol{x_o}) \approx -\eta \sum_{\boldsymbol{x_u} \in \mathcal{D}} \mathcal{A}^t(\boldsymbol{x_o}) \, \mathcal{K}^t(\boldsymbol{x_o}, \boldsymbol{x_u}) \, \mathcal{G}^t(\boldsymbol{x_u}, \boldsymbol{y_u})$$

Updated by gradient of **own label**

Influenced by gradient of other horses

$\pi($ 🐴 $)$

Path on this epoch

$\mathbf{e}_{\mathbf{plane}}$
[**1**, 0, 0]

$\pi^{ES}$: [0.3, 0.6, 0.1]

$\mathbf{e}_{\mathbf{horse}}$
[0, **1**, 0]

- ● epoch start
- ＋ epoch end
- ● Xo update start
- ＋ Xo update end
- — Other Xu update

Easy sample

Medium sample

Update:3696

Hard sample

Hard sample

**Warm up: summary**

$$\Delta \log \pi_{\theta^t}(y|\boldsymbol{x_o}) \approx -\eta \sum_{\boldsymbol{x_u} \in \mathcal{D}} \mathcal{A}^t(\boldsymbol{x_o}) \, \mathcal{K}^t(\boldsymbol{x_o}, \boldsymbol{x_u}) \, \mathcal{G}^t(\boldsymbol{x_u}, \boldsymbol{y_u})$$

✓ Force comes from $\mathcal{G}^t$

✓ Then projected by $\mathcal{K}^t$ and $\mathcal{A}^t$

✓ Finally imposed on $\log \pi(\boldsymbol{x_o})$

✓ $\mathcal{G}^t(\boldsymbol{x_u}, \boldsymbol{y_u})$ evolves with time $t$

# Outline

Chapt. 1 & 2

Chapt. 3 & 4

Chapt. 5

Chapt. 6 & 7



**Intro.**
**PHY & DNN**

**Warm up:**
**Fundamentals of**
**Learning Dynamics**

**Main Application:**
**LLM's fine-tuning**

**Other applications &**
**Theory extensions**

FORCE ANALYSIS OF
A NEURAL NETWORK

*Ren* et al. ICLR 2022
*Ren* et al. ICLR2020
Guo, *Ren* et al. ICLR2024

*Ren* et al. ICLR 2025
*(Outstanding Paper Award)*
Deng, *Ren* et al.. ICML-W 2025

*Ren* et al. ICLR 2023

*Ren* et al. NeurIPS-W 2024
*Ren* et al. NeurIPS 2023
*Ren* et al. ICLR2020

# LEARNING DYNAMICS OF LLM FINETUNING

**Yi Ren**
University of British Columbia
renyi.joshua@gmail.com

**Danica J. Sutherland**
University of British Columbia & Amii
dsuth@cs.ubc.ca

ICLR – 2025
(Outstanding Paper Award)
Chapter 5

# Motivation: unexpected behaviors of SFT

- SFT is good, but many unexpected behaviors:

  ➢ SFT makes the "less preferred responses" more likely

  ➢ SFT exacerbates hallucination

A Closer Look at the Limitations of Instruction Tuning

Sreyan Ghosh         Chandra Kiran Reddy Evuru         Sonal Kumar

Ramaneswaran S         Deepali Aneja         Zeyu Jin

Ramani Duraiswami         Dinesh Manocha

*Siren's Song in the AI Ocean*: A Survey on
Hallucination in Large Language Models

Yue Zhang♠, Yafu Li◇, Leyang Cui♡ , Deng Cai♡, Lemao Liu♡

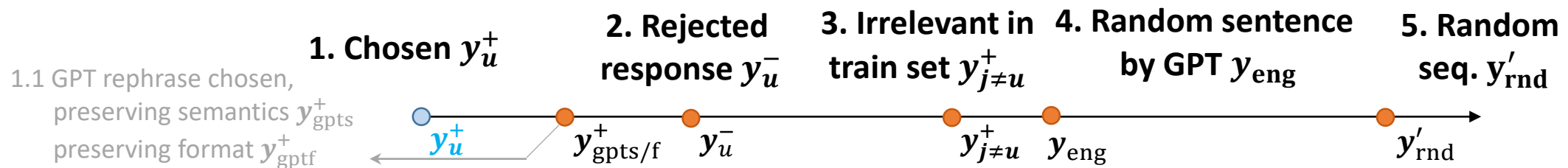# Theory: extend learning dynamics to LLM

$$\chi = [\boldsymbol{x}; \boldsymbol{y}]$$

- After some work, we get:

$$[\Delta \log \pi^t(y|\chi_o)]_m = -\sum_{l=1}^{L} \eta [\mathcal{A}^t(\chi_o)]_m [\mathcal{K}^t(\chi_o, \chi_u)]_{m,l} [\mathcal{G}(\chi_u)]_l + \mathcal{O}(\eta^2)$$

$$\underbrace{\qquad}_{V \times M} \qquad \underbrace{\qquad}_{V \times V \times M} \underbrace{\qquad}_{V \times V \times M \times L} \underbrace{\qquad}_{V \times L}$$

- Check some typical responses (update using $[\boldsymbol{x_u}, \boldsymbol{y_u^+}]$):

**1. Chosen** $y_u^+$    **2. Rejected response** $y_u^-$    **3. Irrelevant in train set** $y_{j \neq u}^+$    **4. Random sentence by GPT** $y_{\text{eng}}$    **5. Random seq.** $y_{\text{rnd}}'$

1.1 GPT rephrase chosen, preserving semantics $y_{\text{gpts}}^+$ preserving format $y_{\text{gptf}}^+$

$y_u^+$    $y_{\text{gpts/f}}^+$    $y_u^-$    $y_{j \neq u}^+$    $y_{\text{eng}}$    $y_{\text{rnd}}'$

**Given question** $x_u$**, our** $y$ **is:**    Valid    Invalid    Ungrammatical

E.g., Antropic-HH, UltraFeedback

$[\mathbf{x_u}, \mathbf{y_u^+}]$ ( 4 , y)    $[\mathbf{x_u}, \mathbf{y_{\text{eng}}}]$ ( 9 , y)    $[\mathbf{x_u}, \mathbf{y_{\text{rnd}}'}]$ $[\mathbf{x_{j \neq u}}, \mathbf{y'}]$ ( 0 , y)
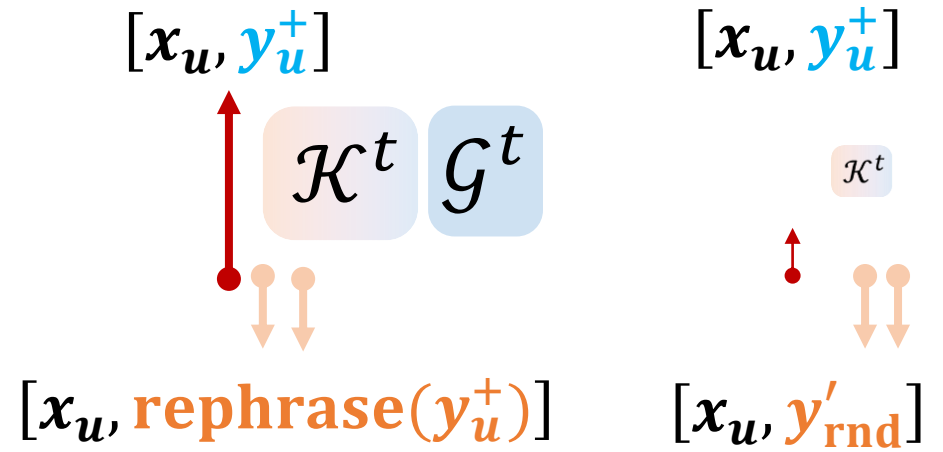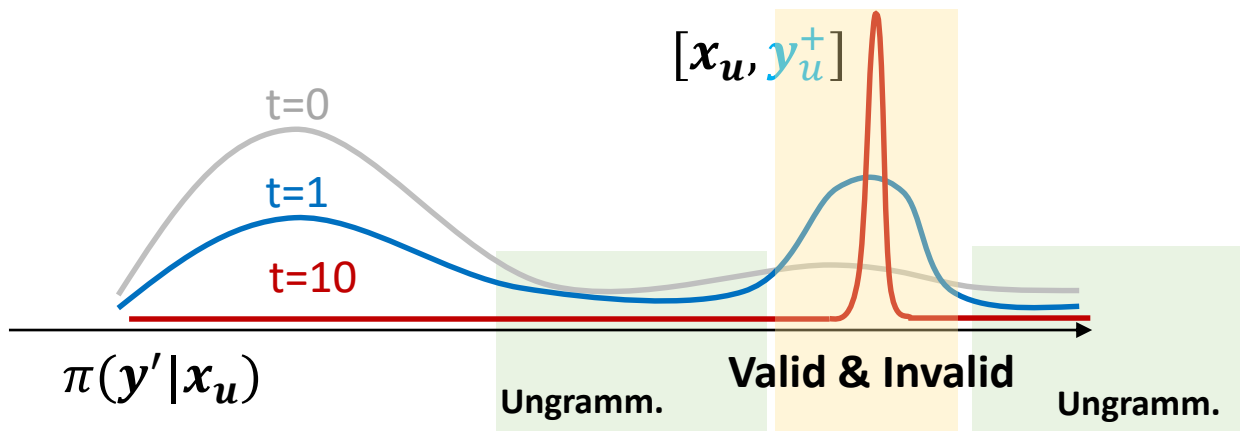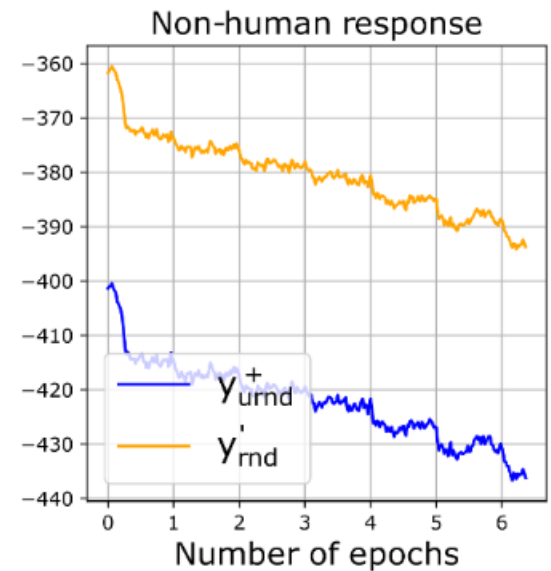
# Application: analyze behaviors in SFT

➢ Why does SFT make the "less prefered answer" more likely?

Because those answers are similar to $[x_u, y_u^+]$

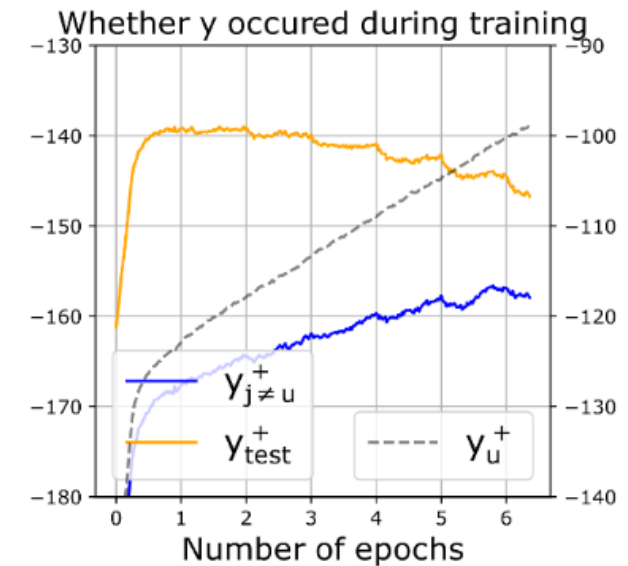$[x_u, y_u^+]$ $\qquad$ $[x_u, y_u^+]$

$\mathcal{K}^t$ $\mathcal{G}^t$ $\qquad$ $\mathcal{K}^t$

$[x_u, \mathbf{rephrase}(y_u^+)]$ $\qquad$ $[x_u, y_{\mathbf{rnd}}']$



$[x_u, y_u^+]$

$t=0$

$t=1$

$t=10$

$\pi(y'|x_u)$

Ungramm.  Valid & Invalid  Ungramm.

Chosen v.s. rejected

Averge log-probability

Number of epochs

$y_u^+$  $y_{gptf}^+$
$y_u^-$  $y_{test}^+$
$y_{gpts}^+$  $y_{eng}$

Non-human response

Number of epochs

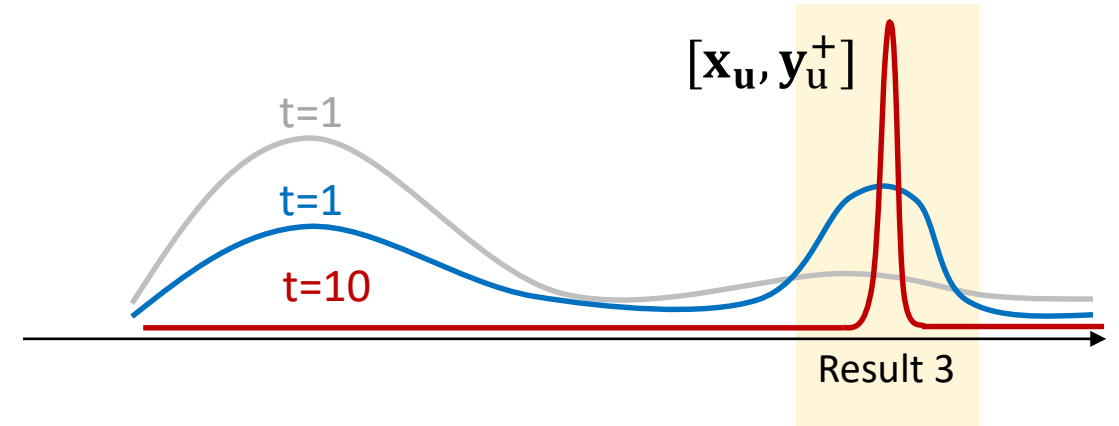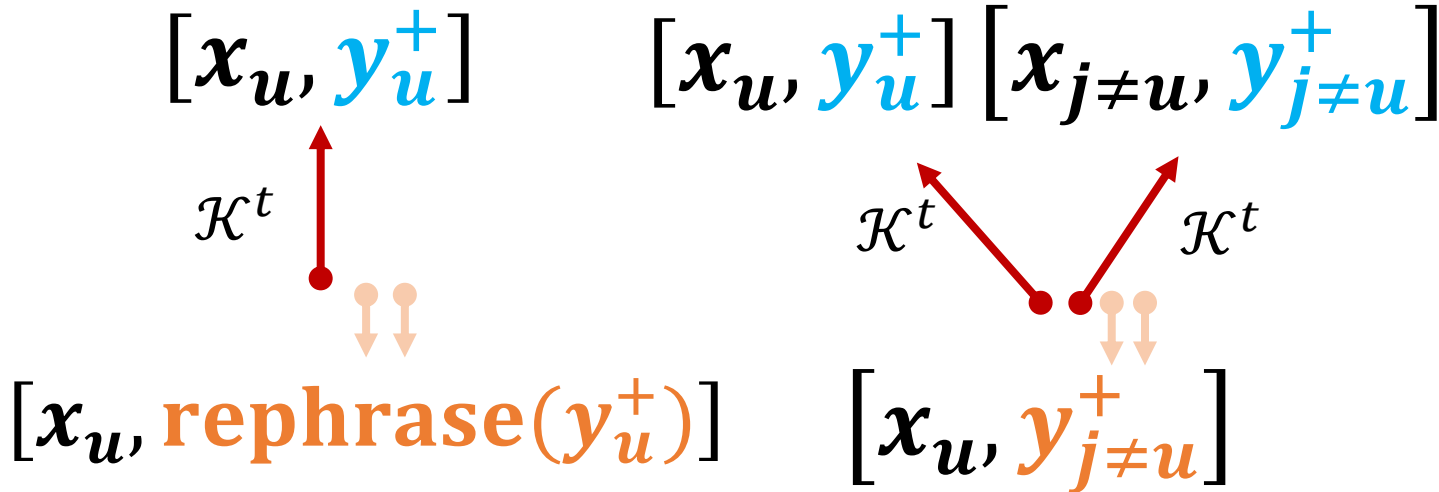$y_{urnd}^+$
$y_{rnd}'$

Valid & Invalid $\qquad$ Ungrammatical

18

# Application: analyze behaviors in SFT

➢ Why does SFT exacerbate hallucination?
   (Specific type of hallucination)

Hallucinated "facts" have MORE "pull-up forces"

$$[x_u, y_u^+]$$

$$\mathcal{K}^t$$

$$[x_u, \text{rephrase}(y_u^+)]$$

$$[x_u, y_u^+] \quad [x_{j \neq u}, y_{j \neq u}^+]$$

$$\mathcal{K}^t \qquad \mathcal{K}^t$$

$$[x_u, y_{j \neq u}^+]$$

$[x_u, y_u^+]$

t=1

t=1

t=10

Result 3



Whether y occured during training

$y_{j \neq u}^+$

$Y_{test}^+$    $y_u^+$

Number of epochs

**Result 3: hallucination!!!**
$[x_u, y_{j \neq u}^+]$ **increases a lot!**
**Using A2 to answer Q1**

19

# More empirical supports: from a famous project

## HALoGEN
### Fantastic LLM Hallucinations and Where to Find Them

Abhilasha Ravichander[1*]    Shrusti Ghela[1†*]    David Wadden[2]    Yejin Choi[13]

https://halogen-hallucinations.github.io/

| Type B | An incorrect fact was in the pretraining data or the fact is taken out of context i.e. the fact appeared within a specific setting in a document in the training data, but when taken in isolation, it loses its original meaning. |
|---|---|

- User Prompt:
*"Write a Python function to calculate the F1 score using scikit-learn."*

- LLM's hallucinated response: $\left[x_u, y_{j \neq u}^+\right]$

```python
from sklearn.metrics import fscore

def calculate_f1(y_true, y_pred):
    return fscore(y_true, y_pred)
```

No `fscore`! Should be `f1_score`

- Where "fscore" comes from:

reddit    $\left[x_{j \neq u}, y_{j \neq u}^+\right]$

… you can calculate **fscore** easily use **sklearn.metrics**, …
… To calculate the **fscore** between two predictions, a straightforward way is to use **sklearn** or pytorch function …

# Motivation: unexpected behaviors in preference tuning

- DPO (or xPO) is good, but more unexpected behaviors:

  ➢ More frequent "repeater" after finetuning

  ➢ DPO makes both $\pi(y+)$ and $\pi(y-)$ decrease

THE CURIOUS CASE OF
NEURAL TEXT DeGENERATION

Ari Holtzman[†‡]    Jan Buys[§†]    Li Du[†]    Maxwell Forbes[†‡]    Yejin Choi[†‡]
[†]Paul G. Allen School of Computer Science & Engineering, University of Washington
[‡]Allen Institute for Artificial Intelligence
[§]Department of Computer Science, University of Cape Town
{ahai,dul2,mbforbes,yejin}@cs.washington.edu, jbuys@cs.uct.ac.za

From $r$ to $Q^*$: Your Language Model is Secretly a Q-Function

Rafael Rafailov*
Stanford University
rafailov@stanford.edu

Joey Hejna*
Stanford University
jhejna@stanford.edu

Ryan Park
Stanford University
rypark@stanford.edu

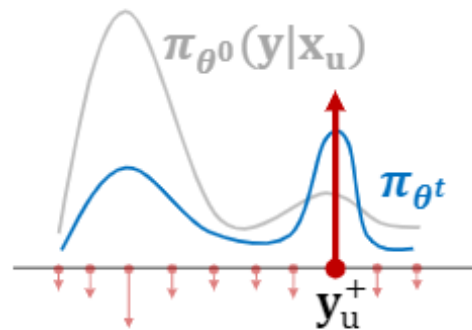Chelsea Finn
Stanford University
cbfinn@stanford.edu

TLDR DPO (β=0.1) implicit reward evolution, no SFT vs. SFT

Legend:
- Chosen rewards, SFT
- Rejected rewards, SFT
- Chosen rewards, no SFT
- Rejected rewards, no SFT

# Theory: extend to DPO, focusing on negative gradient

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(\mathbf{x}_u, \mathbf{y}_u^+, \mathbf{y}_u^-) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta^t}(\mathbf{y}_u^+ \mid \chi_u^+)}{\pi_{\text{ref}}(\mathbf{y}_u^+ \mid \chi_u^+)} \boxed{- \beta \log} \frac{\pi_{\theta^t}(\mathbf{y}_u^- \mid \chi_u^-)}{\pi_{\text{ref}}(\mathbf{y}_u^- \mid \chi_u^-)} \right) \right]$$
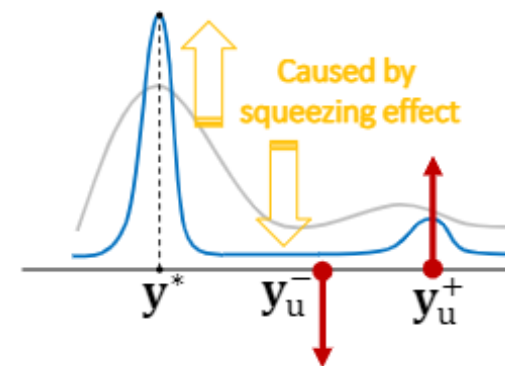
$$[\Delta \log \pi^t(y|\chi_o)]_m \approx -\eta [\mathcal{A}^t(\chi_o)]_m \left( \sum_{l=1}^{L^+} [\mathcal{K}^t(\chi_o, \chi_u^+) \mathcal{G}_{\text{DPO}+}^t]_{m,l} \boxed{- \sum_{l=1}^{L^-} [\mathcal{K}^t(\chi_o, \chi_u^-) \mathcal{G}_{\text{DPO}-}^t]_{m,l}} \right)$$

$$\mathcal{G}_{\text{DPO}+}^t = \beta(1 - \sigma(\cdot))(\pi_{\theta^t}(y|\chi_u^+) - \mathbf{y}_u^+); \quad \mathcal{G}_{\text{DPO}-}^t = \beta(1 - \sigma(\cdot))(\pi_{\theta^t}(y|\chi_u^-) - \mathbf{y}_u^-);$$
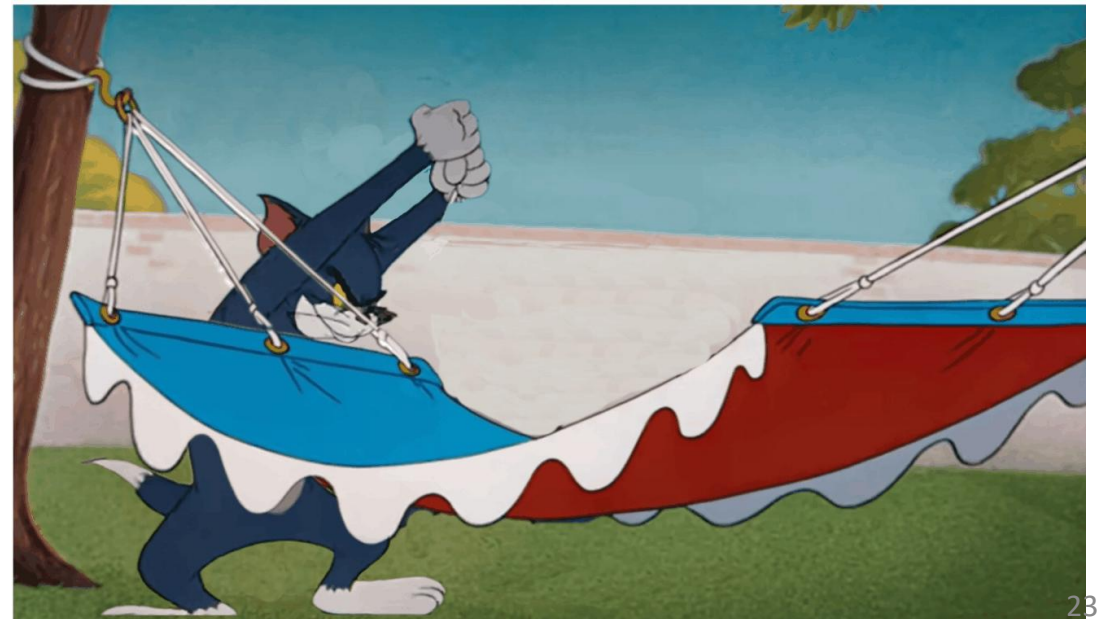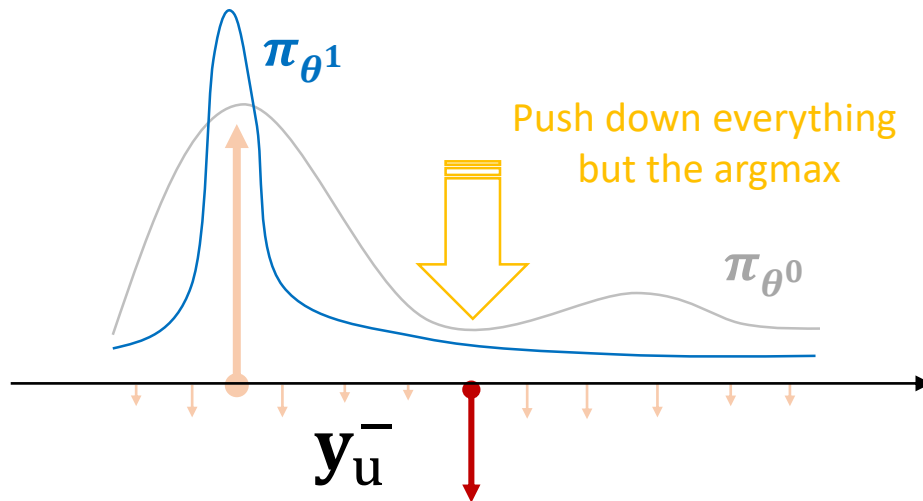


- SFT
- Off-policy DPO, IPO

# Theory: a provable Squeezing Effect !

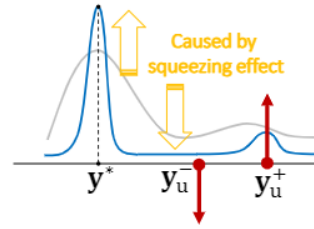- As long as you use Softmax to get probabilies, very likely:

  *Adding **big negative gradient** for an **already unlikely** $y_u^-$*
  *makes weird things happen!*

- (GLOBAL) Almost ALL output probs.$\downarrow\downarrow$
- Except argmax $\uparrow\uparrow$

$$P(y_u^- = 0) = \frac{e^{-10}}{e^{-10} + e^{10} + \cdots}$$

$\pi_{\theta^1}$
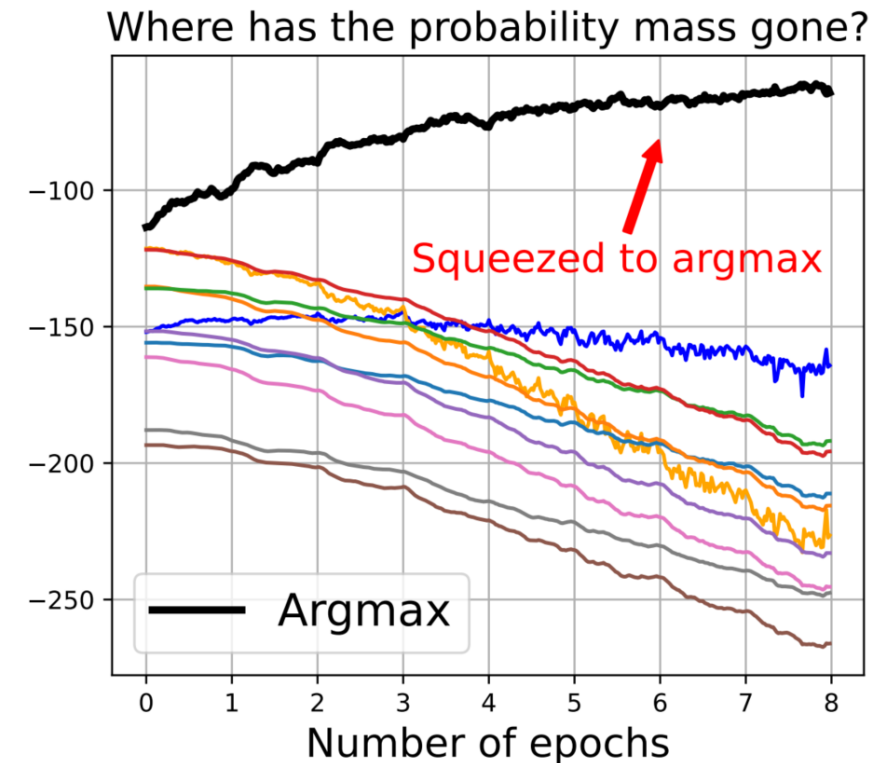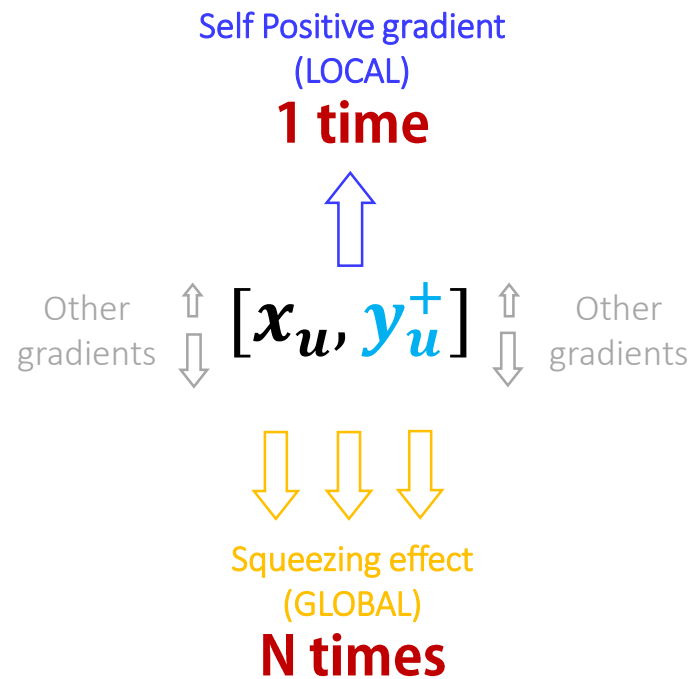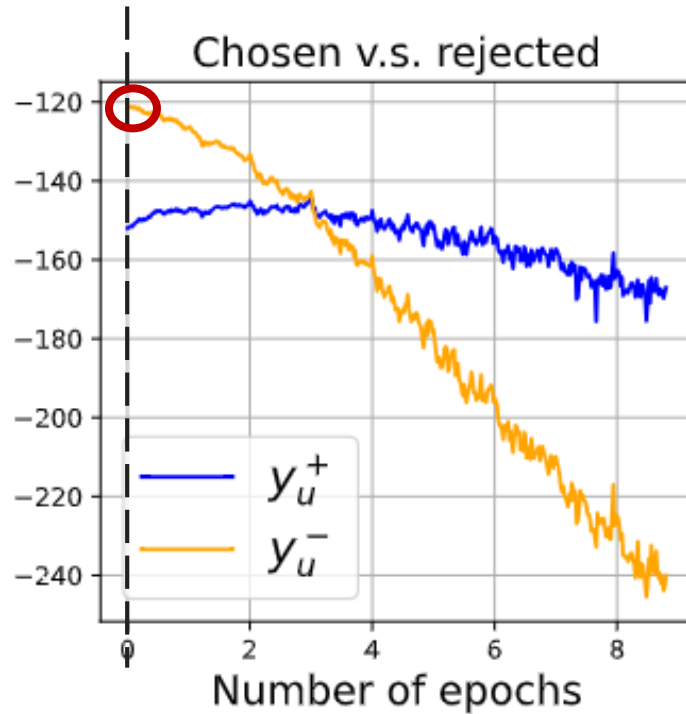
Push down everything
but the argmax

$\pi_{\theta^0}$

$y_u^-$

# Application: analyze off-policy DPO


Caused by squeezing effect

$y^*$  $y_u^-$  $y_u^+$

➢ DPO makes both $\pi$(y+) and $\pi$(y-) decrease
   (Explanation using squeezing effect)

➢ $\pi_\theta(y^*|\chi_u)$ keeps increasing
   (Only self-reinforcing, irrelevant to $\mathcal{D}$)
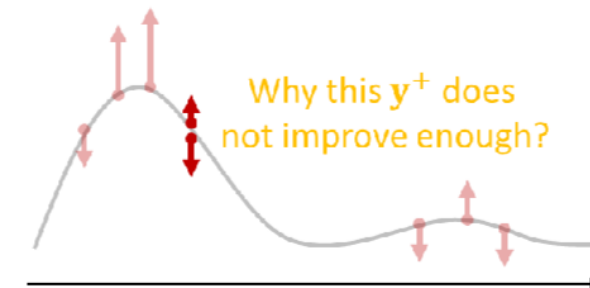
Per-batch (N examples)



Chosen v.s. rejected

Self Positive gradient
(LOCAL)
**1 time**

⇑

Other gradients ⇕  $[\boldsymbol{x_u}, \boldsymbol{y_u^+}]$  ⇕ Other gradients

⇓ ⇓ ⇓

Squeezing effect
(GLOBAL)
**N times**



Where has the probability mass gone?

Squeezed to argmax
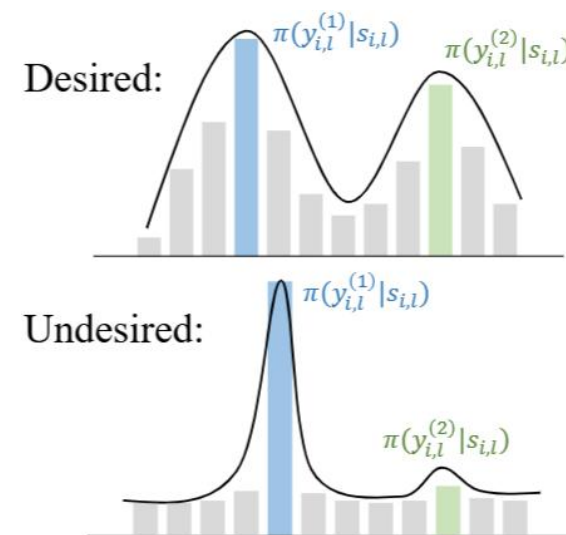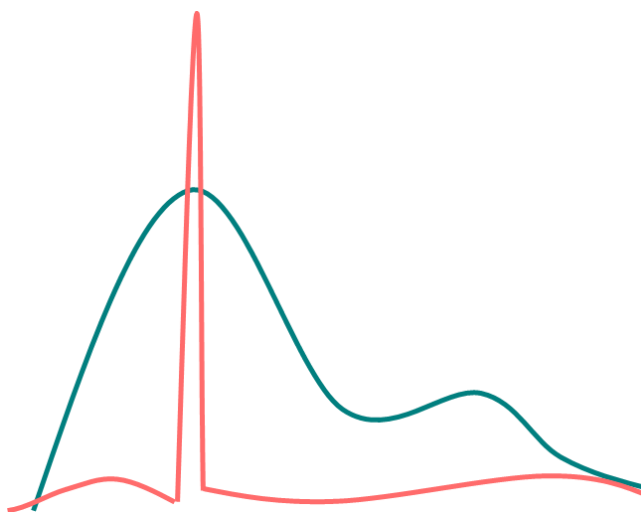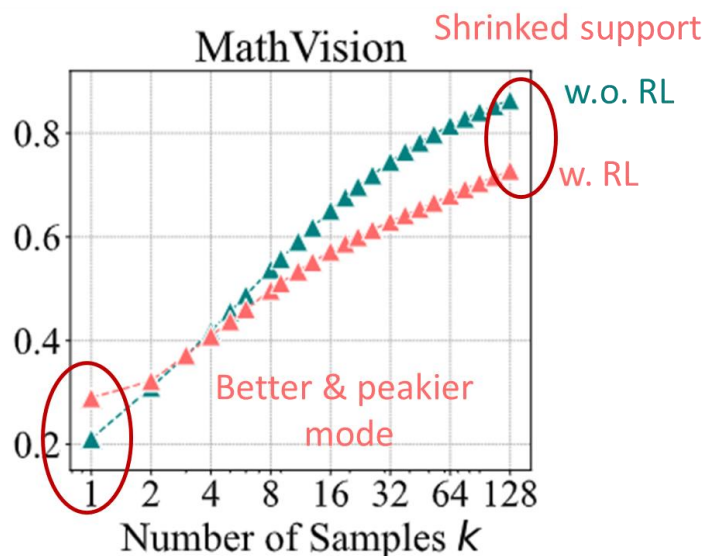
# Application: Improve Exploration in GRPO

➢ Analyze GRPO under the same framework:

$$\mathcal{J}_{\mathrm{GRPO}}(\theta; \gamma_{i,l}) = \frac{1}{G} \sum_{i=1}^{G} \frac{1}{|y_i|} \sum_{l=1}^{|y_i|} [\min(\gamma_{i,l} A_{i,l}, \mathrm{clip}(\gamma_{i,l}, 1-\epsilon, 1+\epsilon) A_{i,l}) - \beta \mathbb{D}_{\mathrm{KL}}(\pi_\theta || \pi_{\mathrm{ref}})]$$

$$\nabla_\theta A_{i,l} \gamma_{i,l} = A_{i,l} \frac{\pi_\theta(y_{i,l}|s_{i,l})}{\pi_{\mathrm{ref}}(y_{i,l}|s_{i,l})} \nabla_\theta \log \pi_\theta(y_{i,l}|s_{i,l}) = \underline{A_{i,l} \cdot \mathrm{sg}(\gamma_{i,l})} \cdot \underline{\nabla_\theta \log \pi_\theta(y_{i,l}|s_{i,l})}$$

Constant Equivalent LR     Same with G-term in SFT and DPO

Why this $\mathbf{y}^+$ does not improve enough?

➢ RLVR hurts exploration ability



Shrinked support
w.o. RL
w. RL
Better & peakier mode

Desired:
$\pi(y_{i,l}^{(1)}|s_{i,l})$   $\pi(y_{i,l}^{(2)}|s_{i,l})$

Undesired:
$\pi(y_{i,l}^{(1)}|s_{i,l})$
$\pi(y_{i,l}^{(2)}|s_{i,l})$

Peng, Ruotian, et al. "SimKO: Simple Pass@ K Policy Optimization." arXiv preprint arXiv:2510.14807 (2025).

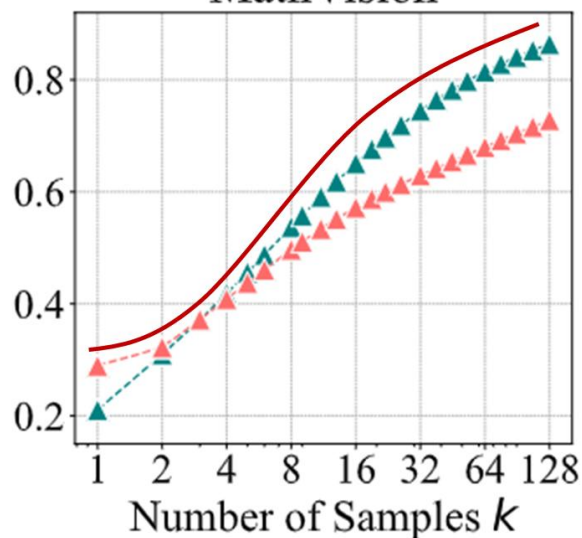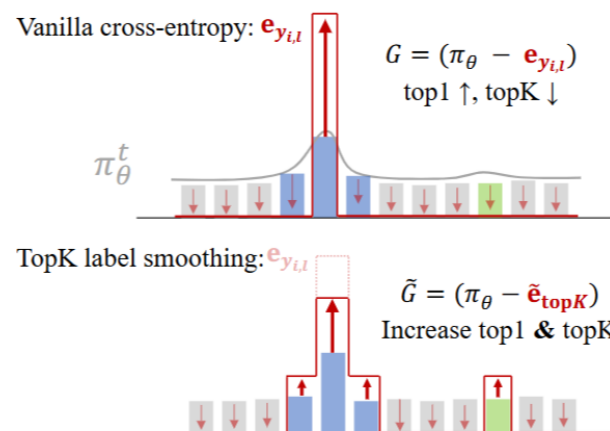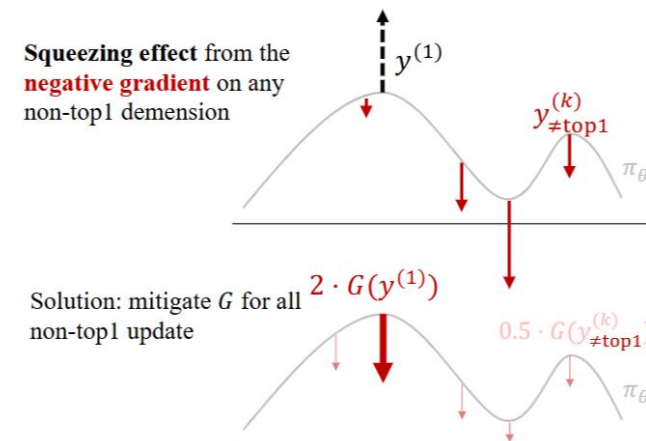# Application: Improve Exploration in GRPO

➢ How to achieve this?

➢ Simple method inspired by learning dynamics

✓ For $A_i > 0$, label smoothing
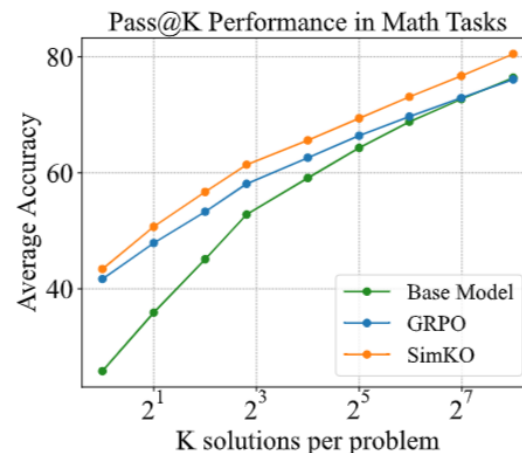
✓ For $A_i < 0$, penalize top1



Vanilla cross-entropy: $\mathbf{e}_{y_{i,l}}$

$G = (\pi_\theta - \mathbf{e}_{y_{i,l}})$
top1 ↑, topK ↓

$\pi_\theta^t$

TopK label smoothing: $\mathbf{e}_{y_{i,l}}$

$\tilde{G} = (\pi_\theta - \tilde{\mathbf{e}}_{topK})$
Increase top1 & topK

**Squeezing effect** from the **negative gradient** on any non-top1 demension

$y^{(1)}$

$y^{(k)}_{\neq top1}$

$\pi_\theta$

Solution: mitigate $G$ for all non-top1 update

$2 \cdot G(y^{(1)})$

$0.5 \cdot G(y^{(k)}_{\neq top1})$

$\pi_\theta$

➢ SimKO results



MathVision



Pass@K Performance in Math Tasks

Pass@K Performance in Logic Tasks

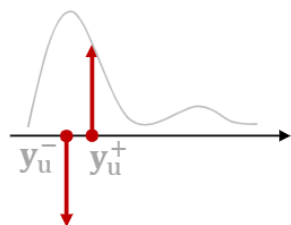Peng, Ruotian, et al. "SimKO: Simple Pass@ K Policy Optimization." arXiv preprint arXiv:2510.14807 (2025).

# LLM Finetuning : summary



- On-policy DPO, IPO
- SPIN
- SPPO
- SLiC
  Triggered when gap is smaller than $\delta$

$y_u^-$ $y_u^+$   $y_u^-$ $y_u^+$   $y_u^-$ $y_u^+$   $y_{ref}$ $y_u^-$ $y_u^+$
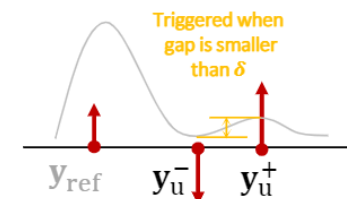
✓ Extension to LLM setting
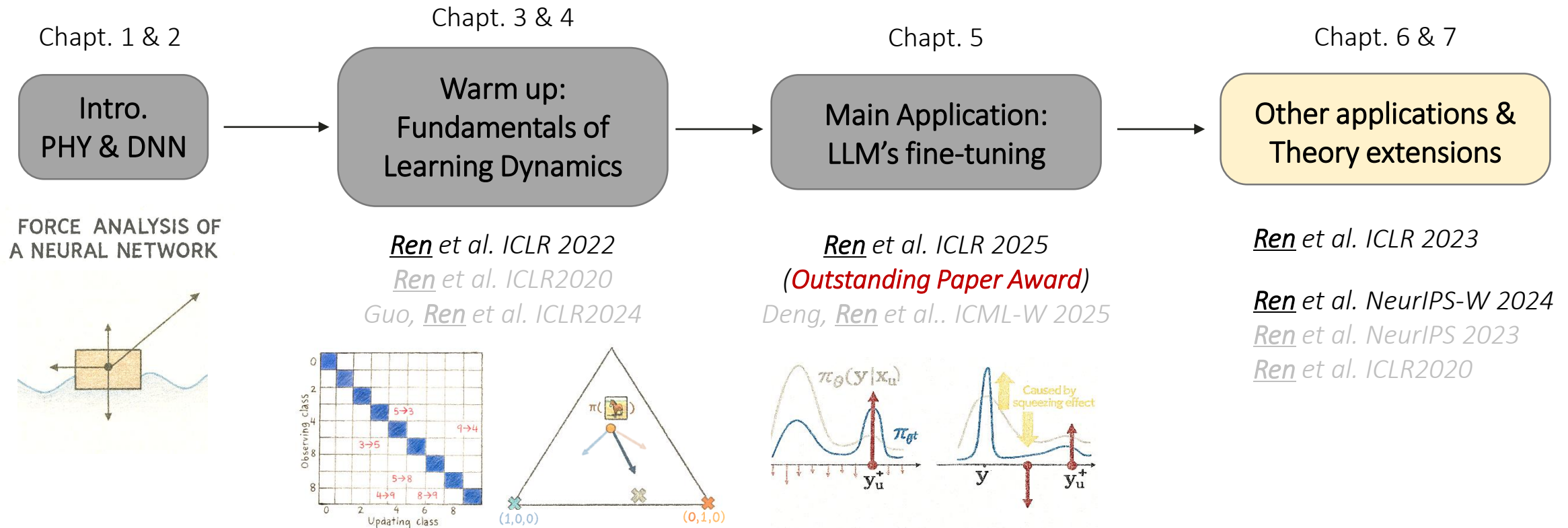(assume relatively stable $\mathcal{K}^t$, more in paper)

✓ Squeezing effect on **negative gradient**
(new findings in the thesis, not coved in ICLR2025 yet!)

✓ Can analyze various methods uniformly
(working on RL-LLMs, using a similar methdology)

# Outline



Chapt. 1 & 2

Chapt. 3 & 4

Chapt. 5

Chapt. 6 & 7

**Intro.**
**PHY & DNN**

**Warm up:**
**Fundamentals of**
**Learning Dynamics**

**Main Application:**
**LLM's fine-tuning**

**Other applications &**
**Theory extensions**

FORCE ANALYSIS OF
A NEURAL NETWORK

*__Ren__ et al. ICLR 2022*
*__Ren__ et al. ICLR2020*
*Guo, __Ren__ et al. ICLR2024*

*__Ren__ et al. ICLR 2025*
*(Outstanding Paper Award)*
*Deng, __Ren__ et al.. ICML-W 2025*

*__Ren__ et al. ICLR 2023*

*__Ren__ et al. NeurIPS-W 2024*
*__Ren__ et al. NeurIPS 2023*
*__Ren__ et al. ICLR2020*

# HOW TO PREPARE YOUR TASK HEAD FOR FINETUNING

**Yi Ren**
University of British Columbia
renyi.joshua@gmail.com

**Shangmin Guo**
University of Edinburgh
s.guo@ed.ac.uk

**Wonho Bae**
University of British Columbia
whbae@cs.ubc.ca

**Danica J. Sutherland**
University of British Columbia & Amii
dsuth@cs.ubc.ca

ICLR – 2023
Chapter 6

# Understanding Simplicity Bias towards Compositional Mappings via Learning Dynamics
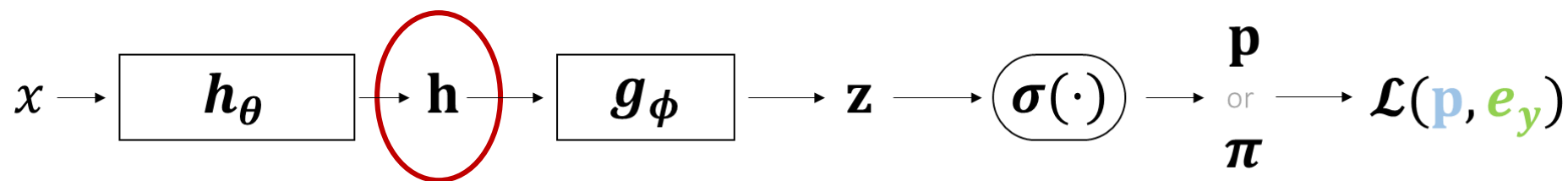
**Yi Ren**
University of British Columbia
renyi.joshua@gmail.com

**Danica J. Sutherland**
University of British Columbia & Amii
dsuth@cs.ubc.ca

NeurIPS Workshop – 2024
Chapter 7

# Chapter 6: understanding general feature adaptation

$$x \rightarrow \boxed{h_\theta} \rightarrow h \rightarrow \boxed{g_\phi} \rightarrow z \rightarrow \sigma(\cdot) \rightarrow \begin{array}{c} p \\ \text{or} \\ \pi \end{array} \rightarrow \mathcal{L}(p, e_y)$$

$$\mathbf{h}_o^{t+1} - \mathbf{h}_o^t = -\eta \frac{1}{N} \sum_{n=1}^{N} \left( \underbrace{\mathcal{K}^t(\mathbf{x}_o, \mathbf{x}_u)}_{\text{slow-change}} \underbrace{(\nabla_\mathbf{h} \mathbf{z}^t(\mathbf{x}_u))^\top}_{\text{direction}} \underbrace{(\mathbf{p}^t(\mathbf{x}_u) - \mathbf{e}_{y_n})}_{\text{energy}} \right) + \mathcal{O}(\eta^2)$$
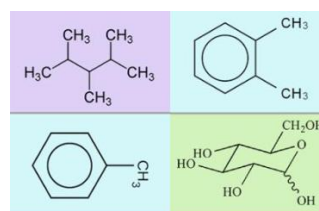
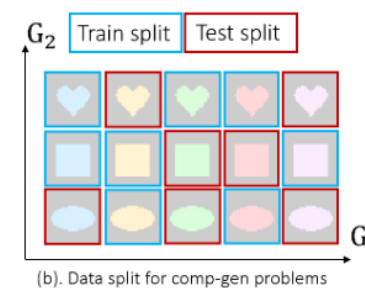- Applying this decomposition to depict how features evolves during training in various deep learning systems:
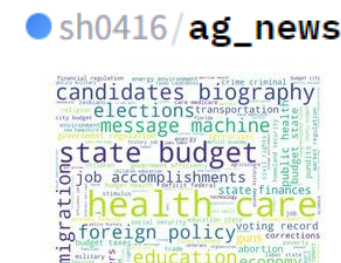
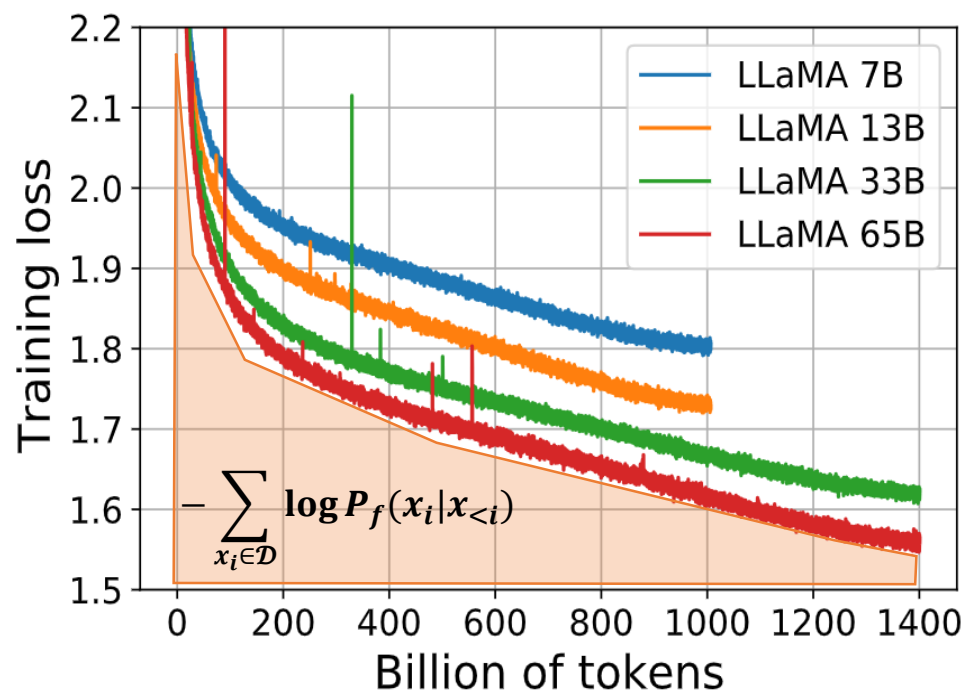| Transfer learning | Image segmantation | Molecular property prediction | Compositional generalization | NLP topic prediction |
|---|---|---|---|---|

_Ren_ et al., ICLR 2023, _Ren_ et al., NeurIPS 2023

# Chapter 7: understanding simplicity bias

- "Compression for AGI" claimed by OpenAI
  (learn faster ⬅➡ better model)



$$-\sum_{x_i \in \mathcal{D}} \log P_f(x_i | x_{<i})$$

Why does this happen spontaneously?

*Ren* et al., NeurIPS-W 2024
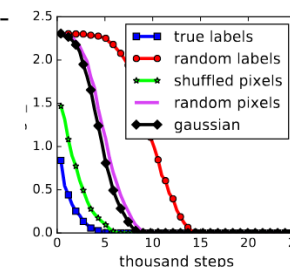
- We provide a novel explanation (in a simple setting):

Good mappings cooperate

Bad mappings contradict

- It can also explain many related phenomena:

UNDERSTANDING DEEP LEARNING REQUIRES RE-
THINKING GENERALIZATION

Clean data learns faster
than noisy labels



A Meta-Transfer Objective for Learning to Disentangle Causal
Mechanisms

Causal data learns faster
than anti-causal

# Thanks for your attention
# Q & A

Yi (Joshua) Ren

renyi.joshua@gmail.com

| Intro. PHY & DNN | | Warm up: Fundamentals of Learning Dynamics | | Main Application: LLM's fine-tuning | | Other applications & Theory extensions |

FORCE ANALYSIS OF
A NEURAL NETWORK

*Ren* et al. ICLR 2022
*Ren et al. ICLR2020*
*Guo, Ren et al. ICLR2024*

*Ren* et al. ICLR 2025
*(Outstanding Paper Award)*
*Deng, Ren et al.. ICML-W 2025*

*Ren* et al. ICLR 2023

*Ren* et al. NeurIPS-W 2024
*Ren et al. NeurIPS 2023*
*Ren et al. ICLR2020*

Good cooperate
Bad contradict