



清华大学 电子工程系

Department of Electronic Engineering, Tsinghua University

# Model Compression Towards Efficient Deep Learning Inference

Xuefei Ning

Department of Electronic Engineering, Tsinghua University

2023.03

foxdoraame@gmail.com

Lab Leader: Prof. Yu Wang [yu-wang@tsinghua.edu.cn](mailto:yu-wang@tsinghua.edu.cn)





## 目录 Contents

- 1 Why does efficient DL inference matters?
- 2 How to generally accelerate DL inference?
- 3 Research on specific applications
  - Low-level vision
  - Image generation
  - 3D perception



## 目录 Contents

- 1 Why does efficient DL inference matters?
- 2 How to generally accelerate DL inference?
- 3 Research on specific applications
  - Low-level vision
  - Image generation
  - 3D perception

# Goal of Artificial Intelligence



**Goal: To serve humanity and improve human life**



**Fei-Fei Li**  
**Co-Director of HAI**

—"ImageNet Project"

"Artificial Intelligence has the potential to help us realize our shared dream of a better future **for all of humanity**. ..., our vision is led by our commitment to studying, guiding and developing **human-centered AI technologies and applications**."

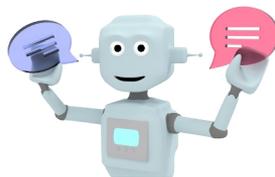
**Interaction with Human**



Human



AI System



- Understand **human's instructions**
- Improve their **sensory experiences**

**Interaction with World**



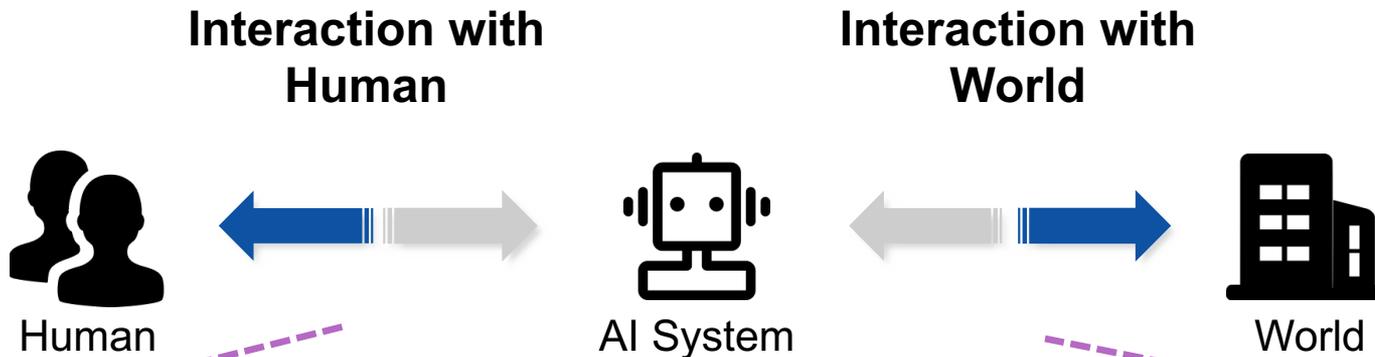
World



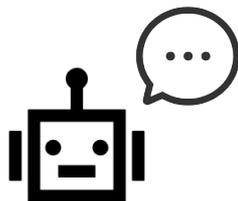
- Understand the **physical world**
- Make decisions to **influence the world**



# Pathway to Improve the Interaction



## Generative tasks



Dialogue Generation



Image Generation



Code Generation



3D Assets Generation

# DL Trend: Data and Task



## Big data, multimodal, generative tasks

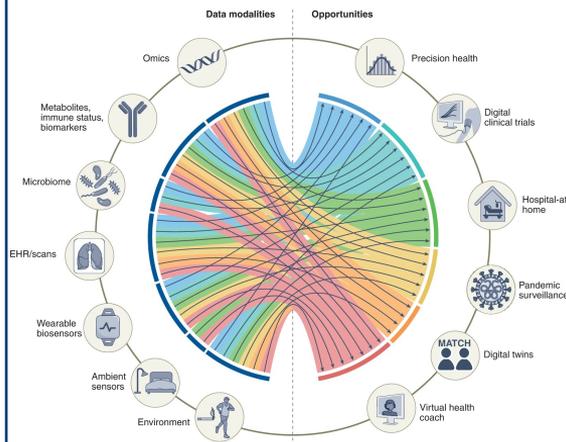


### Big Data



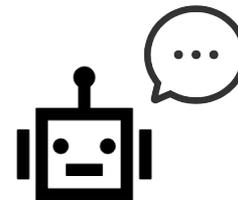
90% data in last 2 years [1]

### Multimodal



multimodal data in healthcare [2]

### Generative Tasks



Text/Code Symbol Generation



Image/3D assets generation

generative tasks are attracting attentions

[1] MIT 6.S191 Course: <https://www.youtube.com/watch?v=QZxcTZj0L-M>

[2] Acosta, J.N., Falcone, G.J., Rajpurkar, P. et al. Multimodal biomedical AI. *Nat Med* **28**, 1773–1784 (2022). <https://doi.org/10.1038/s41591-022-01981-2>

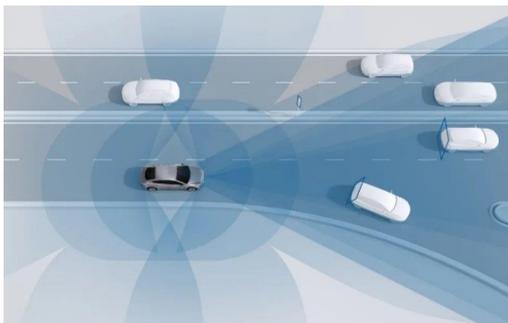
# DL Trend: Data and Task



## Big data, multimodal, generative tasks



### Big Data

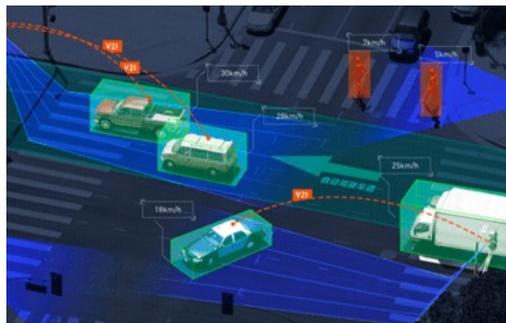


L4

Autonomous Driving

GB/Day/Car → TB/Day/Car

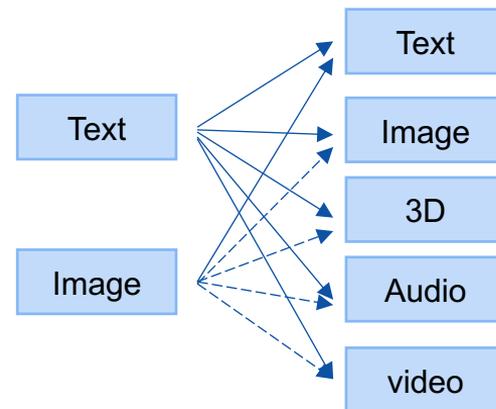
### Multimodal



Intelligent perception of  
vehicle road collaboration

5-6 types, 30-40 sensors

### Generative Tasks



Enormous generative  
applications

Source: EqualOcean Intelligence 2021, Computing Power Driven Vehicles - 2021 Research Report on the Development of Computing Power of China's Intelligent Vehicles



# Various Application Scenario

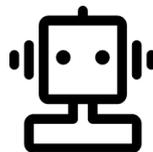


Interaction with Human

Interaction with World



Human



AI System



World



Sensor  
Wearable Device



Mobile Phone  
IoT Device



Smart City  
Auto-driving Car

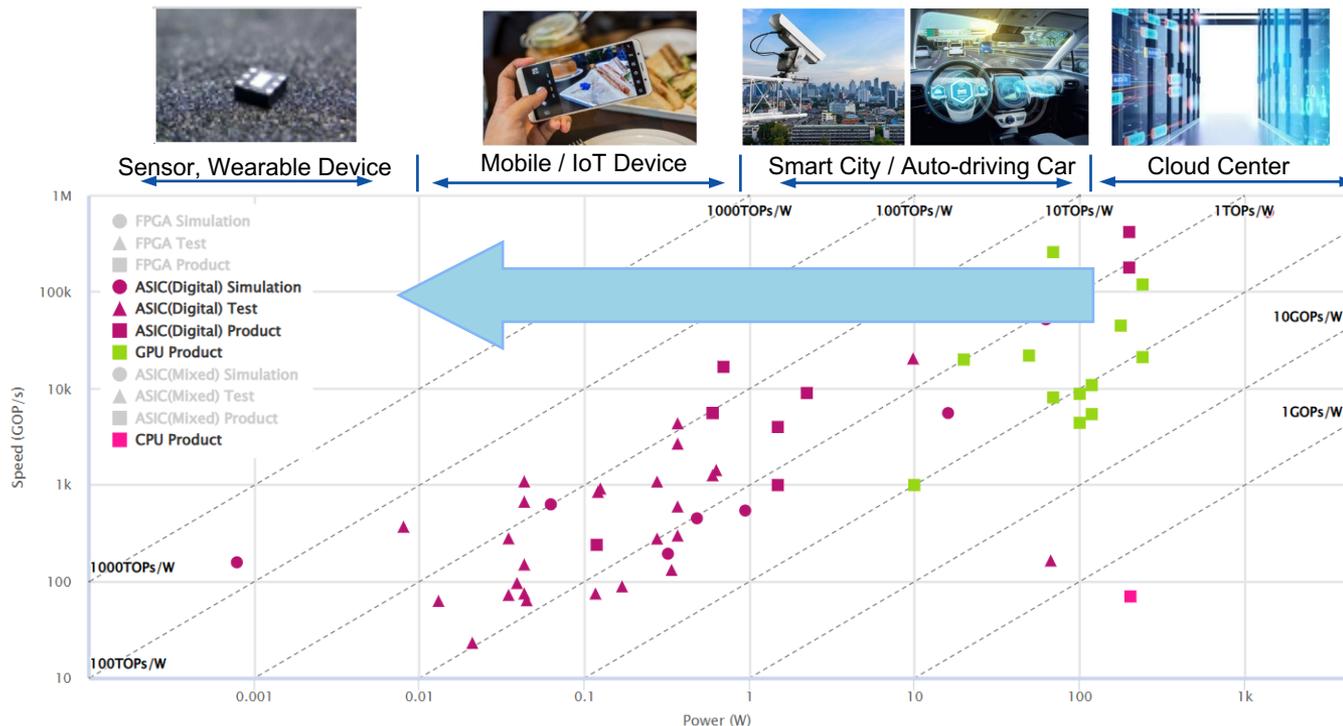


Smart City  
Auto-driving Car

# DL Trend: Computing Power



From cloud center to tiny edge device



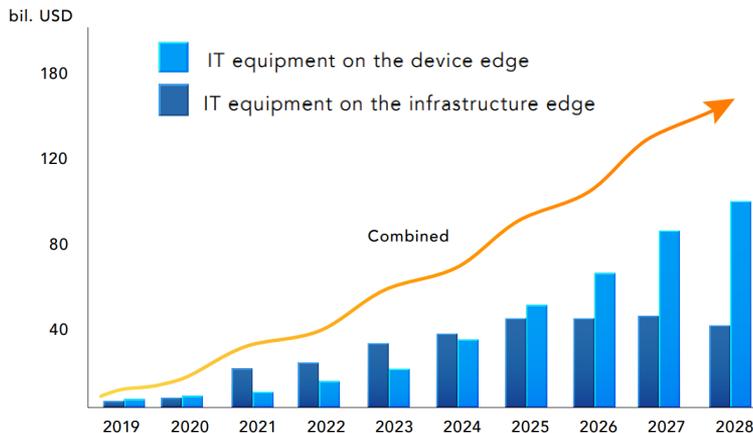
# DL Trend: Computing Power



## From cloud center to tiny edge device

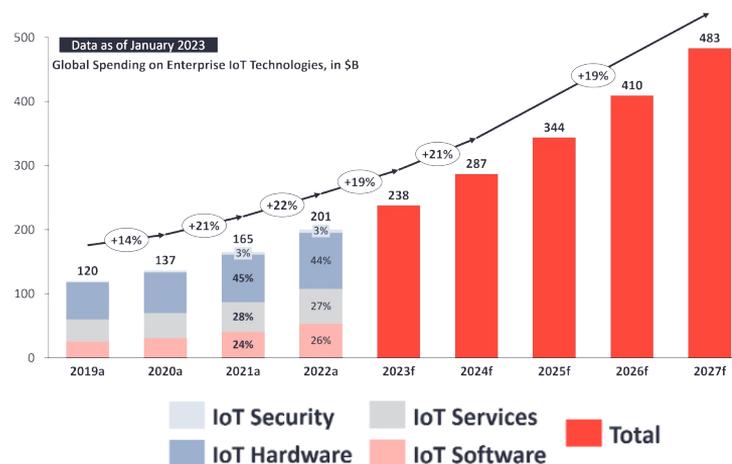


### Global Annual CAPEX Spend on Edge [1]



**Edge: 88% v.s. Cloud: 12%**

### Enterprise IoT market 2019-2027 [2]



**IoT market grew 21.5% in 2022**

[1] State of the Edge 2021: A Market and Ecosystem Report for Edge Computing

[2] <https://iot-analytics.com/iot-market-size/>

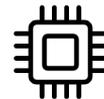
# DL Trend: Summary



Data &  
Task



Model &  
Algorithm



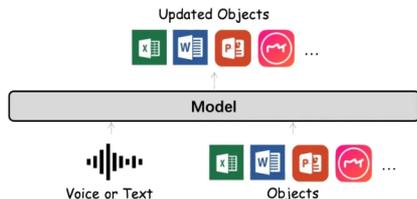
Computing  
Power



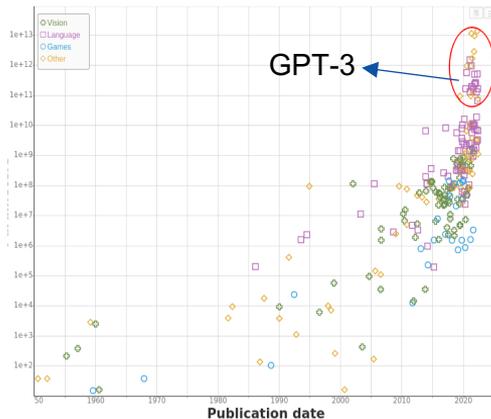
Big data



Multimodal



Generative tasks



Large model



Cloud AI  
Data centers  
Expensive  
Connection required  
Privacy issue



Mobile AI  
Smartphones  
Accessible  
Process locally



Tiny AI  
IoT Devices/  
Microcontrollers  
Cheap, small, low-power  
Rapid growth



Tiny edge device

# Application Challenges

- **Gap** caused by the large model and the tiny edge device

## Real-time Requirement

on Constrained Platform



Video Conference

~30FPS

Our research goal is to **reduce the inference latency of the model to satisfy the application requirement** (e.g., latency on certain platform)

Outdoor Detection

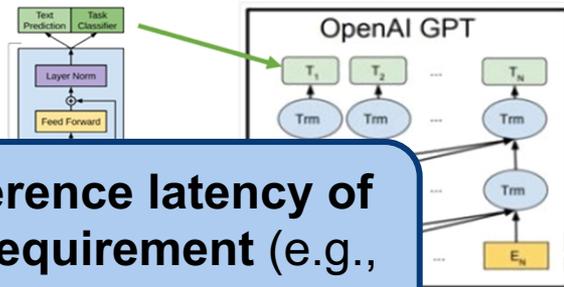
10~25Hz

Auto-driving Car

(10~100 TOPS)

## High Latency

caused by Large Model Size



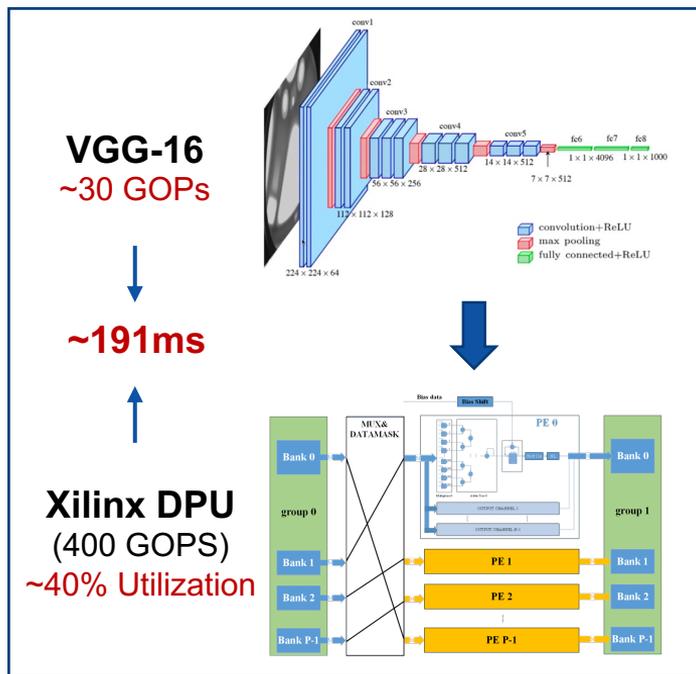
<1FPS

NVIDIA A100  
(19.5 TFLOPS)

E.g., stable diffusion: 1.45B, 5.7TOPs, **2s/img**  
GPT-3: 175B, 664TOPs, **32s/seq**

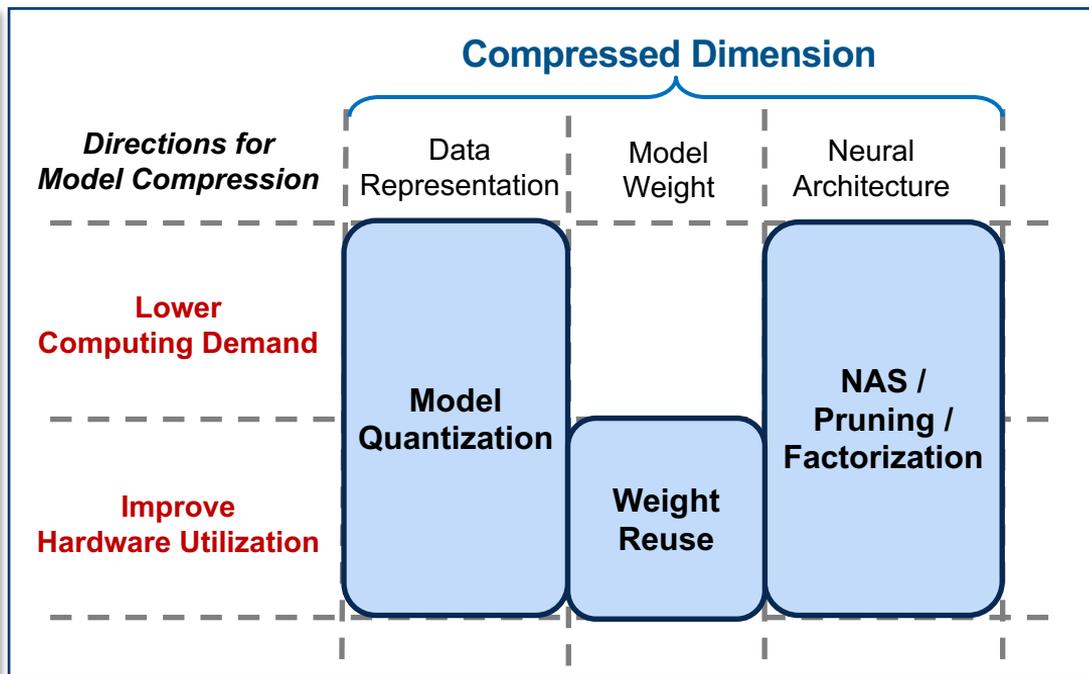
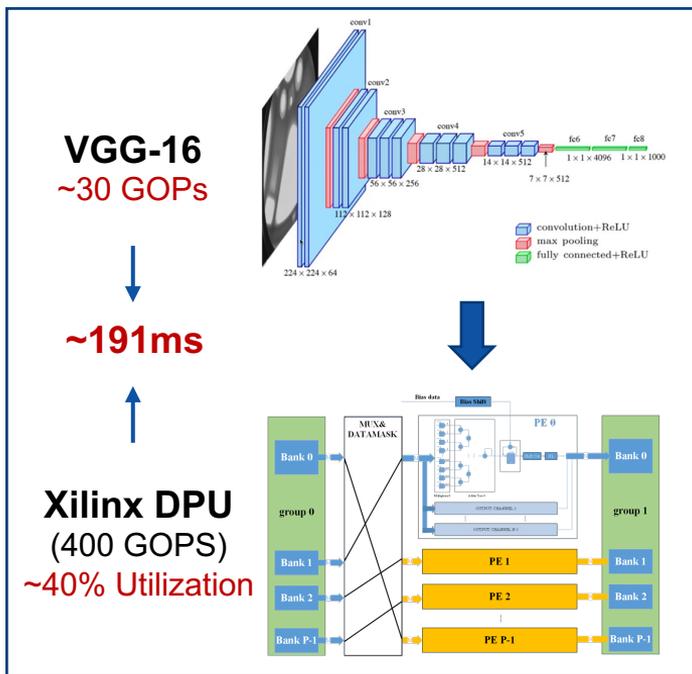
# Direction to Mitigate the Gap

$$\text{Latency} = \frac{\text{Computing Demand}}{\text{Computing Performance} * \text{Hardware Utilization}}$$



# Direction to Mitigate the Gap

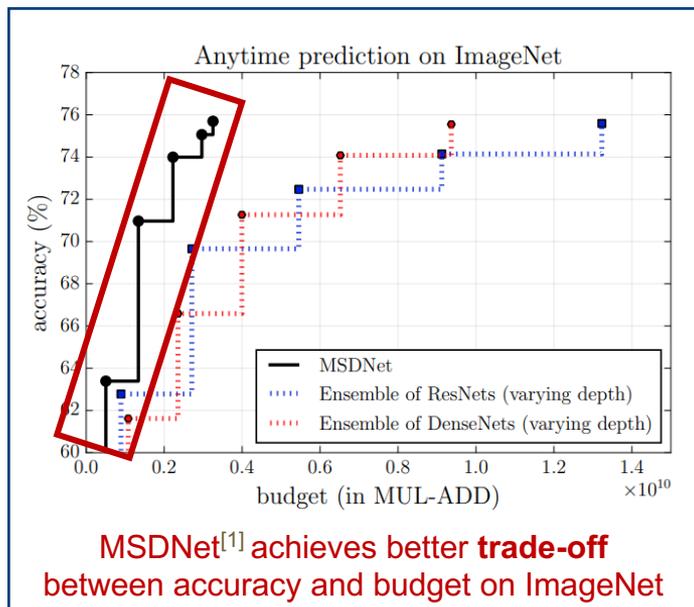
$$\downarrow \text{Latency} = \frac{\text{Computing Demand} \downarrow}{\text{Computing Performance} * \text{Hardware Utilization} \uparrow}$$



# Direction to Mitigate the Gap



## *System Performance: Task Performance, Latency*

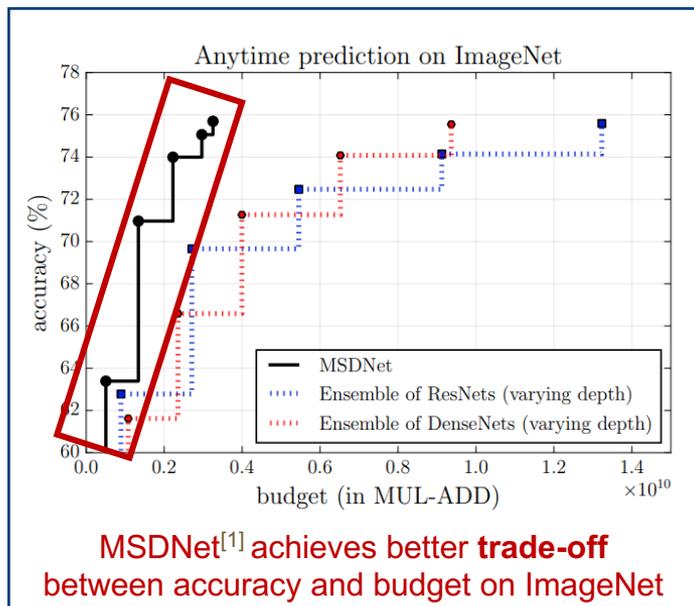


[1] Huang et al., Multi-Scale Dense Convolutional Networks for Efficient Prediction, ICLR'18 oral

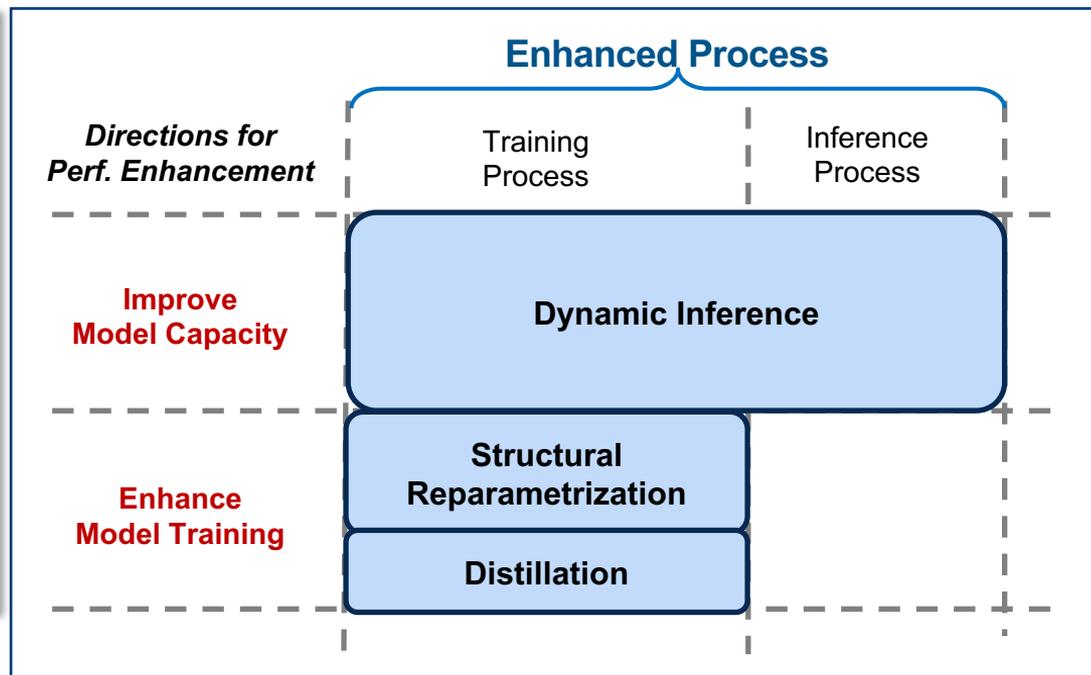
# Direction to Mitigate the Gap



↑ *System Performance: Task Performance*, ↑ *Latency*



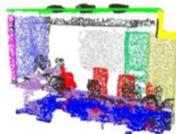
[1] Huang et al., Multi-Scale Dense Convolutional Networks for Efficient Prediction, ICLR'18 oral



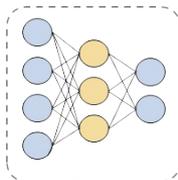
# Automated Model Compression

- Varying tasks and hardware may need different neural network
- Need to consider multiple objectives or constraints in the meantime

## Varying Tasks



"A little girl sits in a plastic swing set."



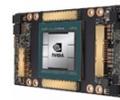
Neural Network

Expensive/Suboptimal



Purely Manual Model Compression

## Varying Hardware



A100 GPU



Versal ACAP



TPU v3



NVIDIA Xavier GPU



Xilinx ZU11 FPGA



TI TDA4 ASIC

We need automated model compression!



## 目录 Contents

- 1 Why does efficient AI inference matters?
- 2 How to generally accelerate AI inference?
- 3 Research on specific applications
  - Low-level vision
  - Image generation
  - 3D perception

# Model Compression

- Basic Problem Definition

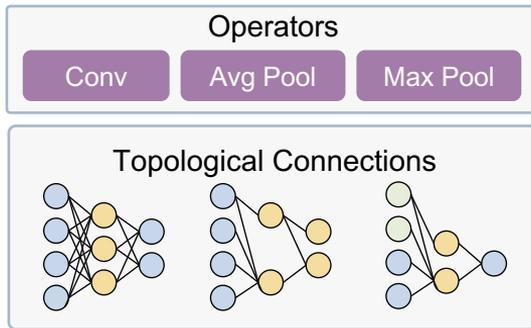
- Neural network architecture  $a$ ; Search space  $S$ ; Network parameters  $w$
- Valid task perf  $R_{val}$ ; Train task loss  $L_{train}$ ; Complexity function  $F$ ; Budget  $B$

$$\underset{a \in S}{\text{Maximize}} \quad R(w^*(a), a) \quad \rightarrow \quad \text{Targets: task perf, perf loss, ...}$$

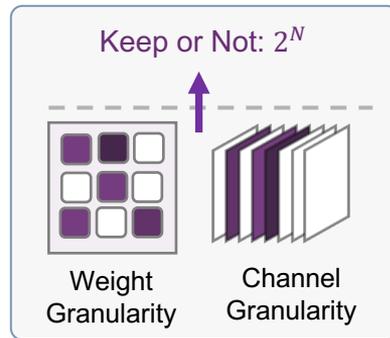
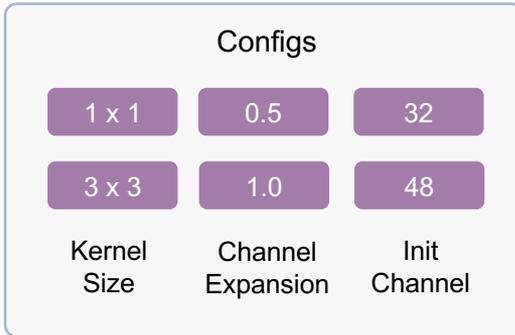
$$\text{s.t.} \quad w^*(a) = \underset{w}{\text{argmin}} L_{train}(w, a)$$

$$F(a) \leq B \quad \rightarrow \quad \text{Constraints: latency, throughput, energy, ...}$$

- Optimize in the **Network Architecture** dimension: **NAS / Pruning**



**NAS Search Space**



**Pruning Space**

# Model Compression

- Basic Problem Definition

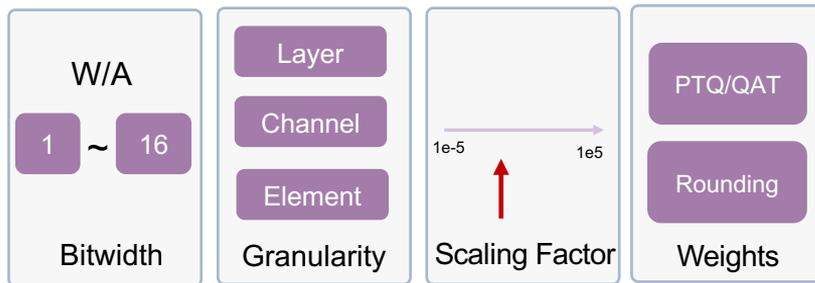
- Neural network architecture  $a$ ; Search space  $S$ ; Network parameters  $w$
- Valid task perf  $R_{val}$ ; Train task loss  $L_{train}$ ; Complexity function  $F$ ; Budget  $B$

$$\underset{a \in S}{\text{Maximize}} \mathbf{R}(w^*(a), a) \rightarrow \text{Targets: task perf, perf loss, ...}$$

$$\text{s.t. } w^*(a) = \underset{w}{\text{argmin}} L_{train}(w, a)$$

$$\mathbf{F}(a) \leq \mathbf{B} \rightarrow \text{Constraints: latency, throughput, energy, ...}$$

- Optimize in the **Data Representation** dimension: **Quantization**

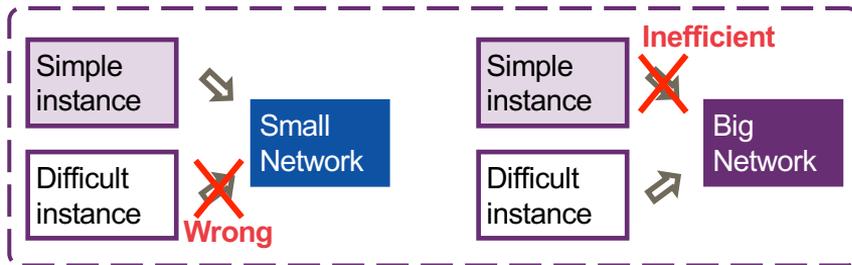


**Quantization Space**

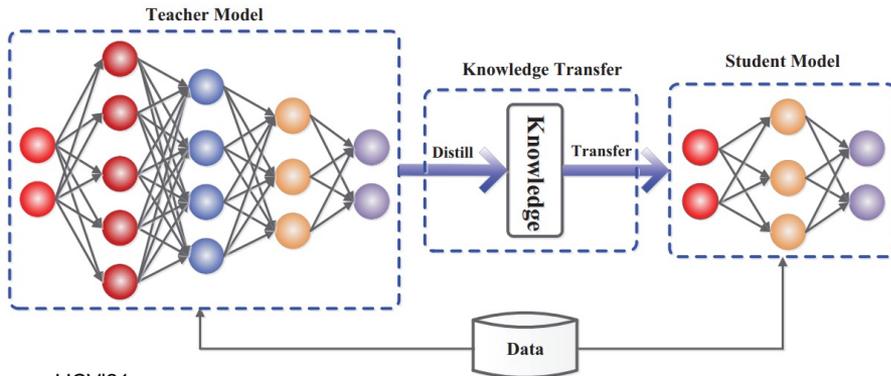
# Enhanced Techniques

- Optimize the **model capacity**: **Dynamic Inference**

Better trade-off between accuracy & efficiency

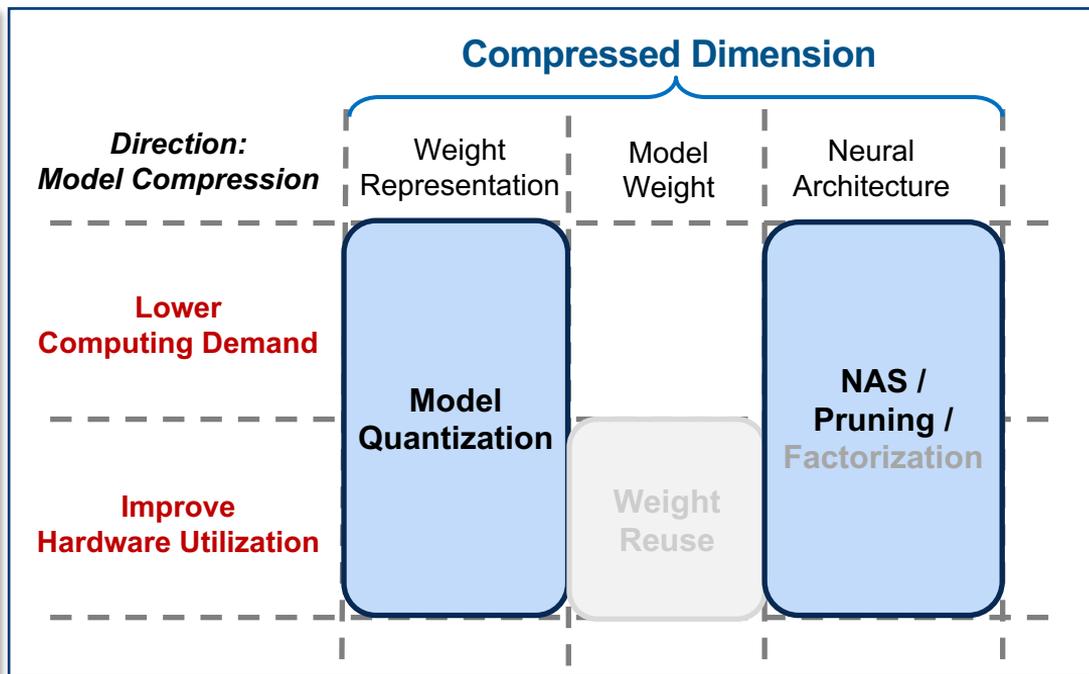
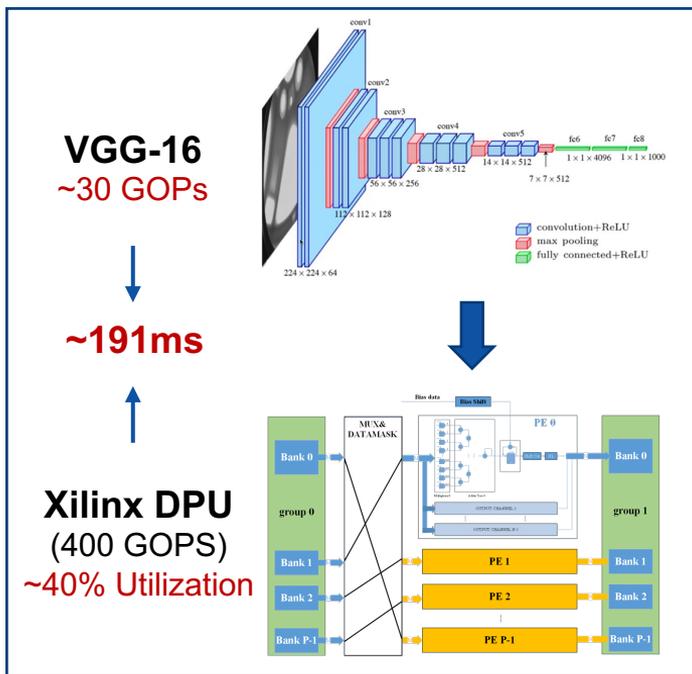


- Enhance the **model training**: **Distillation**



# (Review) Direction to Mitigate the Gap

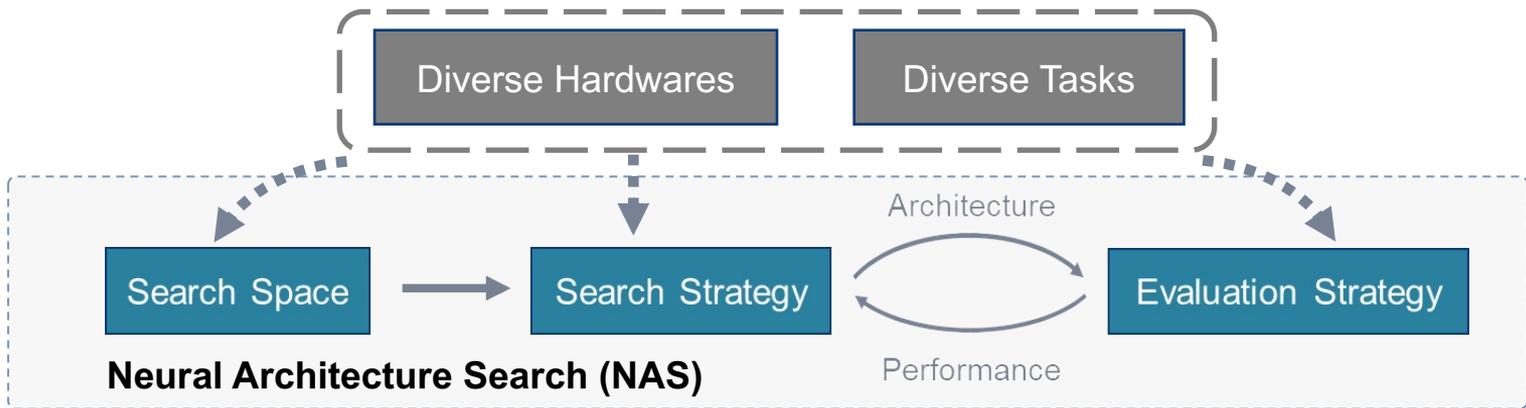
$$\downarrow \text{Latency} = \frac{\text{Computing Demand} \downarrow}{\text{Computing Performace} * \text{Hardware Utilization} \uparrow}$$





# Neural Architecture Search

- Definition: Given the **hardware-related constraints** (e.g., latency, FLOPs), search for the **network architecture** to satisfy the constraints and maximize the accuracy
- Main Components:
  - **Search Space**: Define the searchable configurations, the space of candidate architectures
  - **Evaluation Strategy**: How to evaluate the architecture
  - **Search Strategy**: How to explore the search space



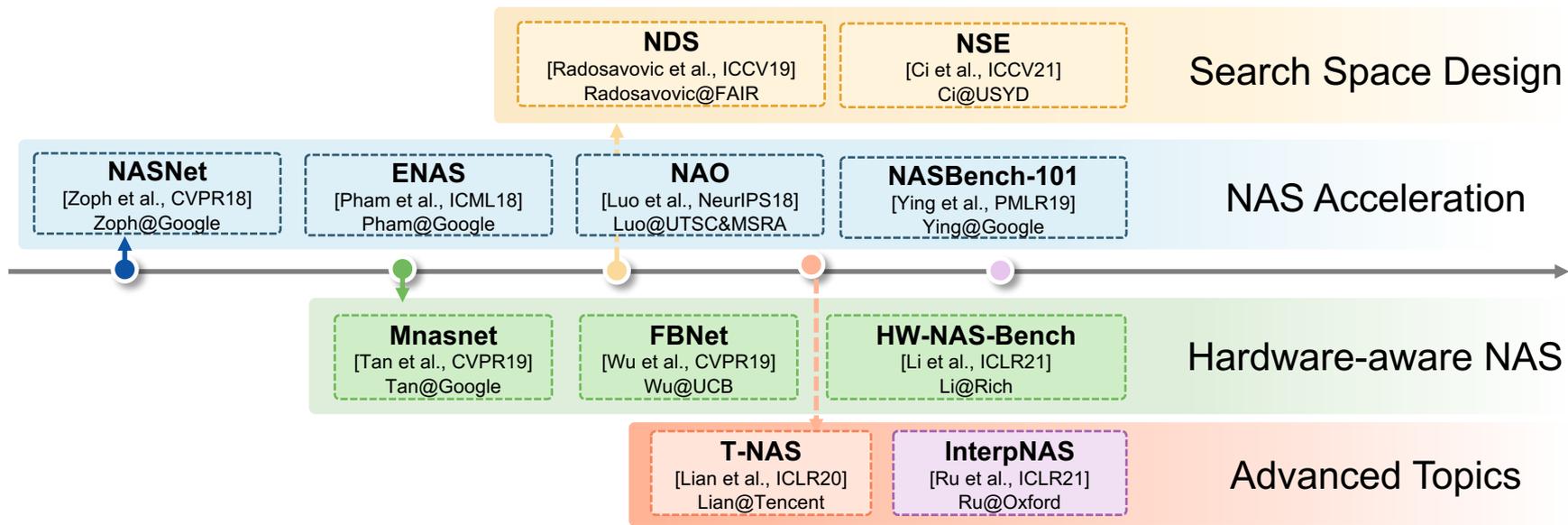
Refer to our websites <https://sites.google.com/view/nas-nicsefc> for more introduction on NAS.

# NAS Development Trends



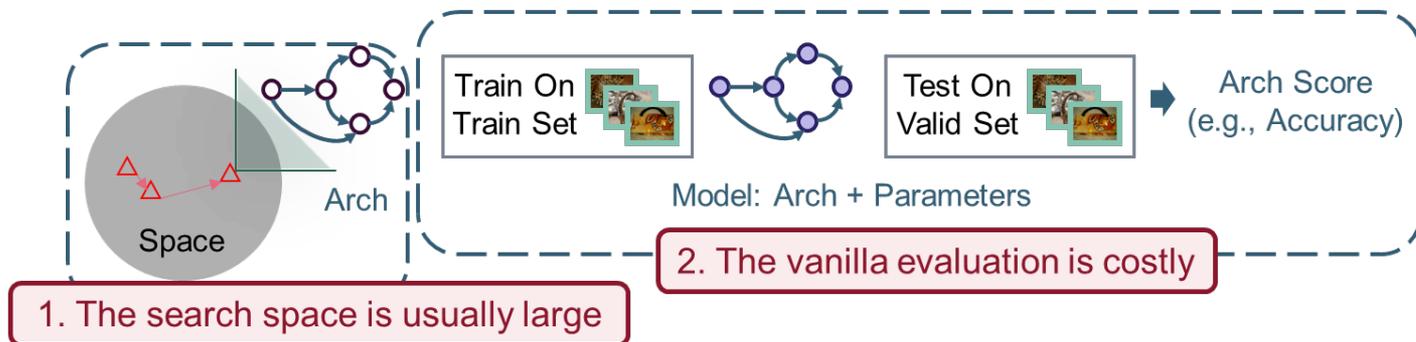
## • Trends

- **Improve the tradeoff between NAS efficiency and performance**
- Take the hardware system into consideration → co-design with system configurations
- Manually design search space → Automatically design
- ...



# Challenge

- Key Challenge: NAS **process is slow**
  - Large search space: Need to evaluate many architectures to explore sufficiently
  - Costly vanilla evaluation: The evaluation of each architecture is costly



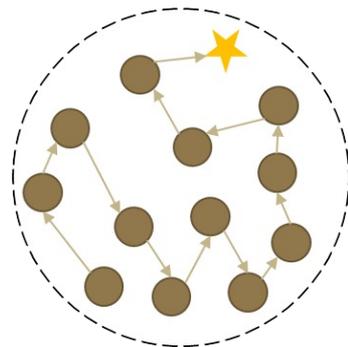
Work	Search space size	#Evaluated architectures	#Train epochs	#GPU day
NASRL [Zoph et al., 2017]	$3.9 \times 10^{73}$	12.8k	50	800x28d=22.4k
NASNet [Zoph et al. 2018]	$5.2 \times 10^{33}$	20k	20	500x4d=2k
AmoebaNet-A [Real et al. 2019]	$3.1 \times 10^{28}$	20k	25	450x7d=3.2k

Zoph et al., "Neural architecture search with reinforcement learning", ICLR'17.

Zoph et al., "Learning Transferable Architectures for Scalable Image Recognition", CVPR'18.

Real et al., "Regularized Evolution for Image Classifier Architecture Search", AAAI'19.

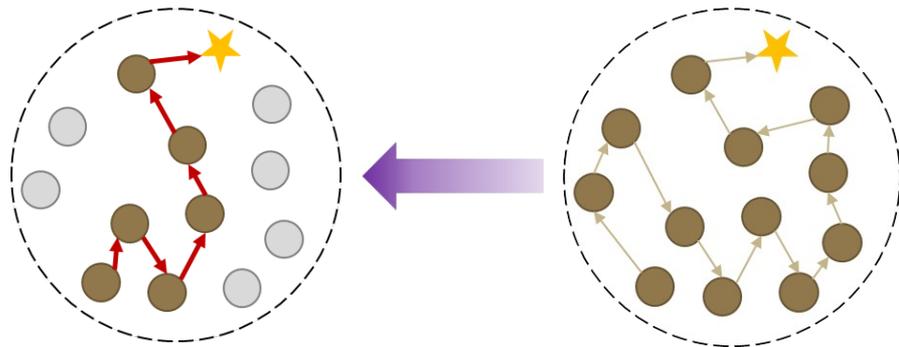
# Direction towards Efficient NAS



# Direction towards Efficient NAS

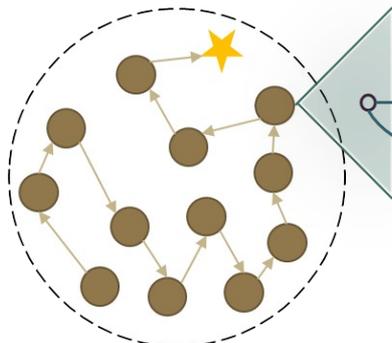
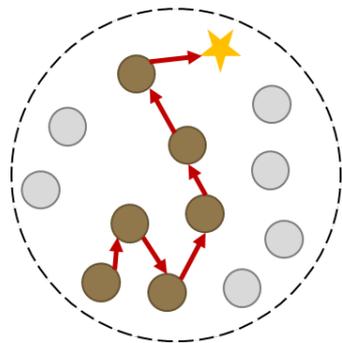


1. Improve the sample efficiency  
of the search strategy

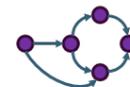


# Direction towards Efficient NAS

1. Improve the sample efficiency  
of the search strategy



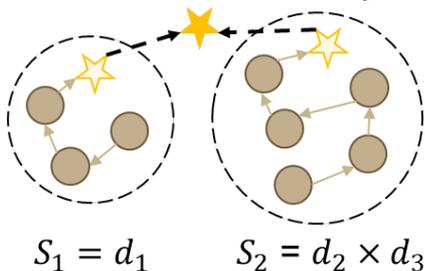
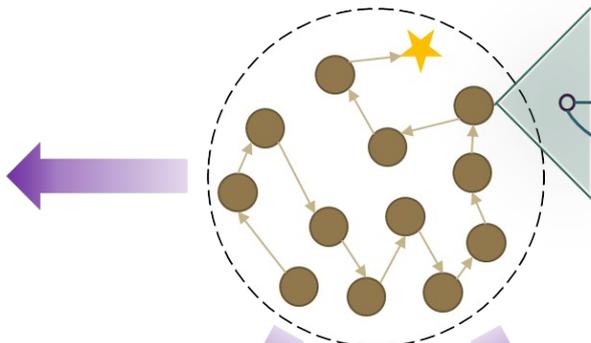
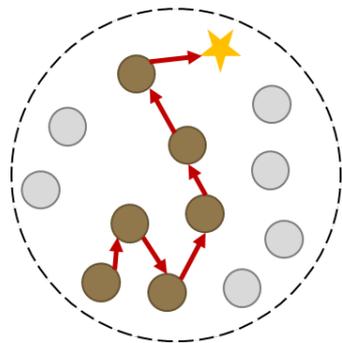
Standalone Training  
on Train Set



  
Low Efficiency

# Direction towards Efficient NAS

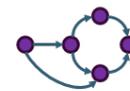
1. Improve the sample efficiency of the search strategy



2. Accelerate the evaluation strategy

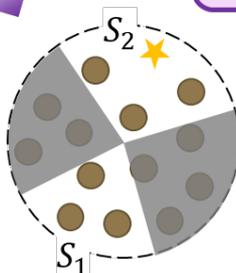


Standalone Training on Train Set



Low Efficiency 

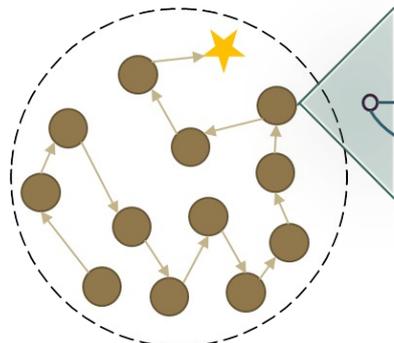
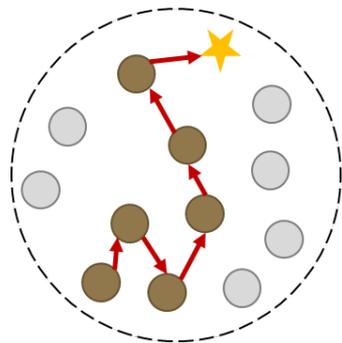
3. Factorize or partition the search space to reduce space complexity.



# Our Solution towards Efficient NAS

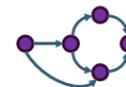


1. Improve the sample efficiency  
of the search strategy



2. Accelerate the evaluation  
strategy

Standalone Training  
on Train Set



Low Efficiency



**Architecture  
Modeling**

Neural Architecture  
Description (DAG)



**Architecture  
Encoder**

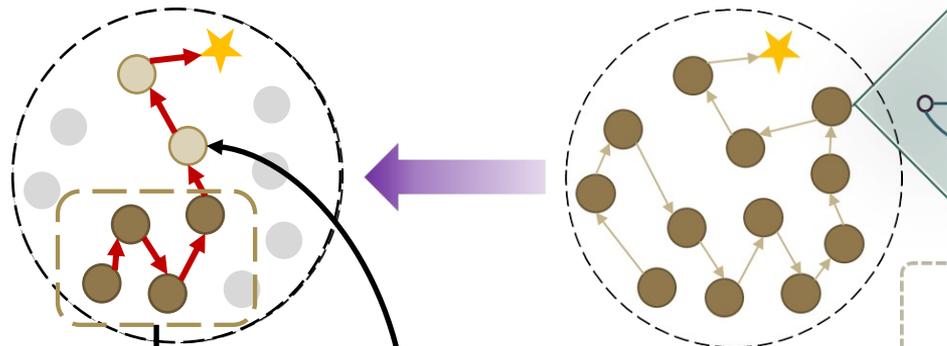
Continuous  
Embedding



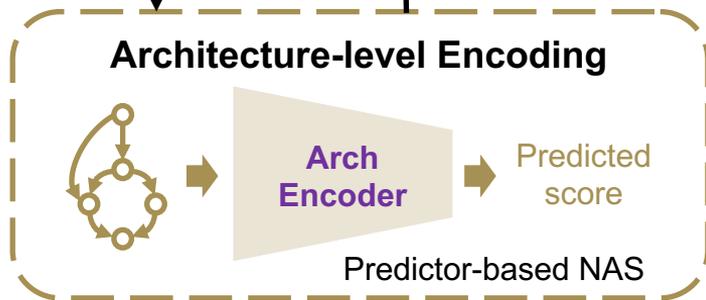
# Our Solution towards Efficient NAS

1. Improve the sample efficiency of the search strategy

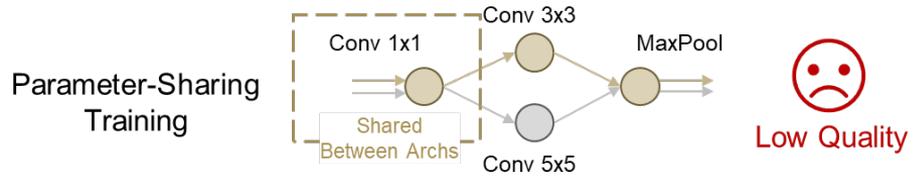
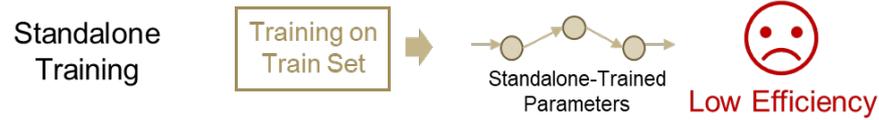
2. Accelerate the evaluation strategy



Standalone Training on Train Set

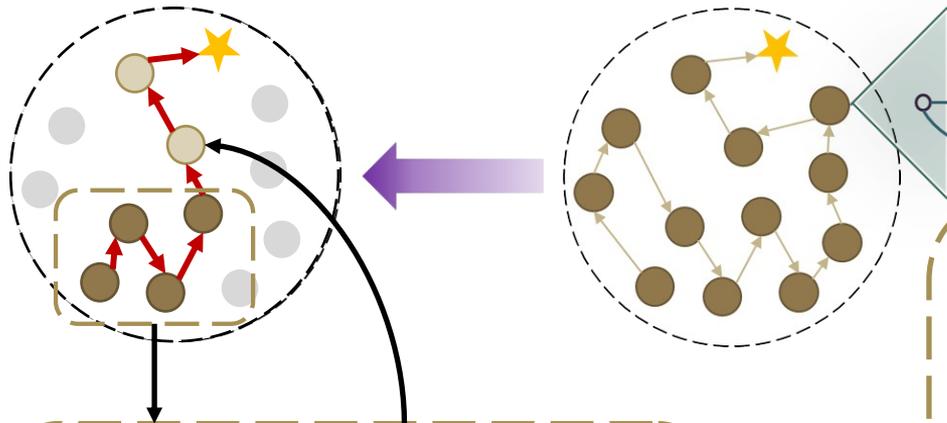


Evaluation: How to get the parameters

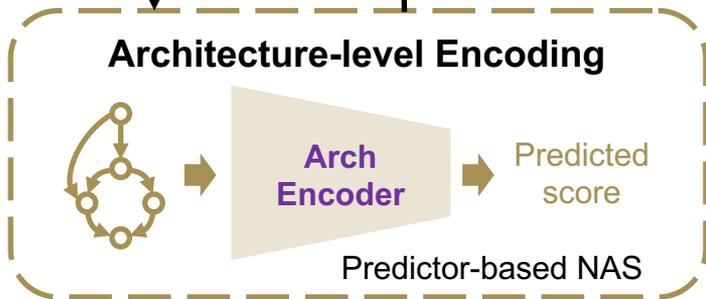
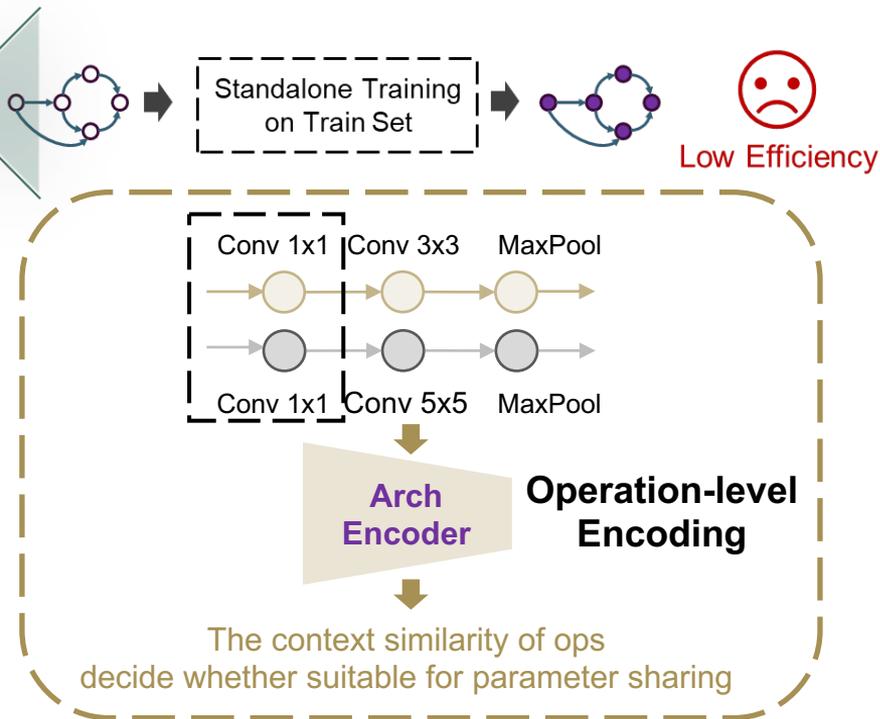


# Our Solution towards Efficient NAS

1. Improve the sample efficiency of the search strategy

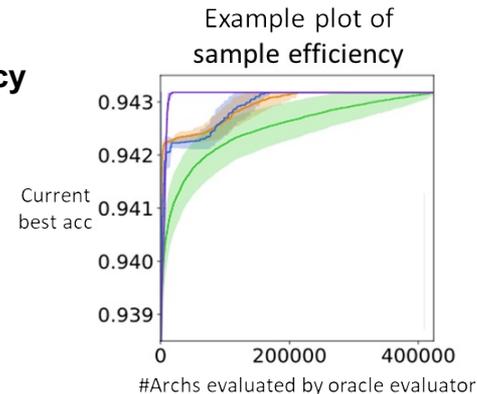
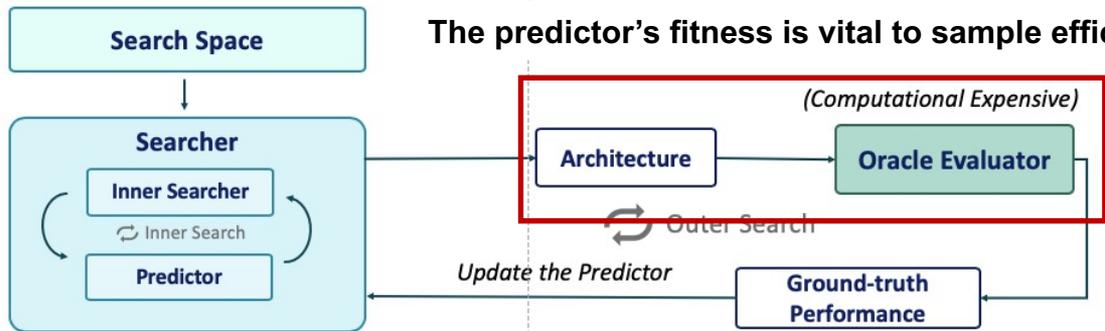


2. Accelerate the evaluation strategy

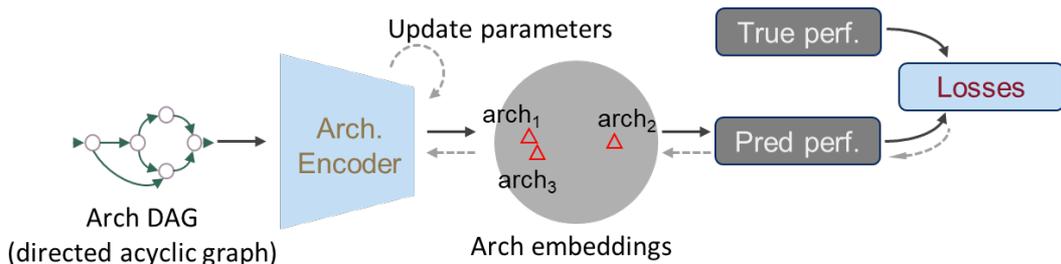


# Enhance the Predictor-based NAS

- Predictor-based NAS



- Typical construction of the predictor

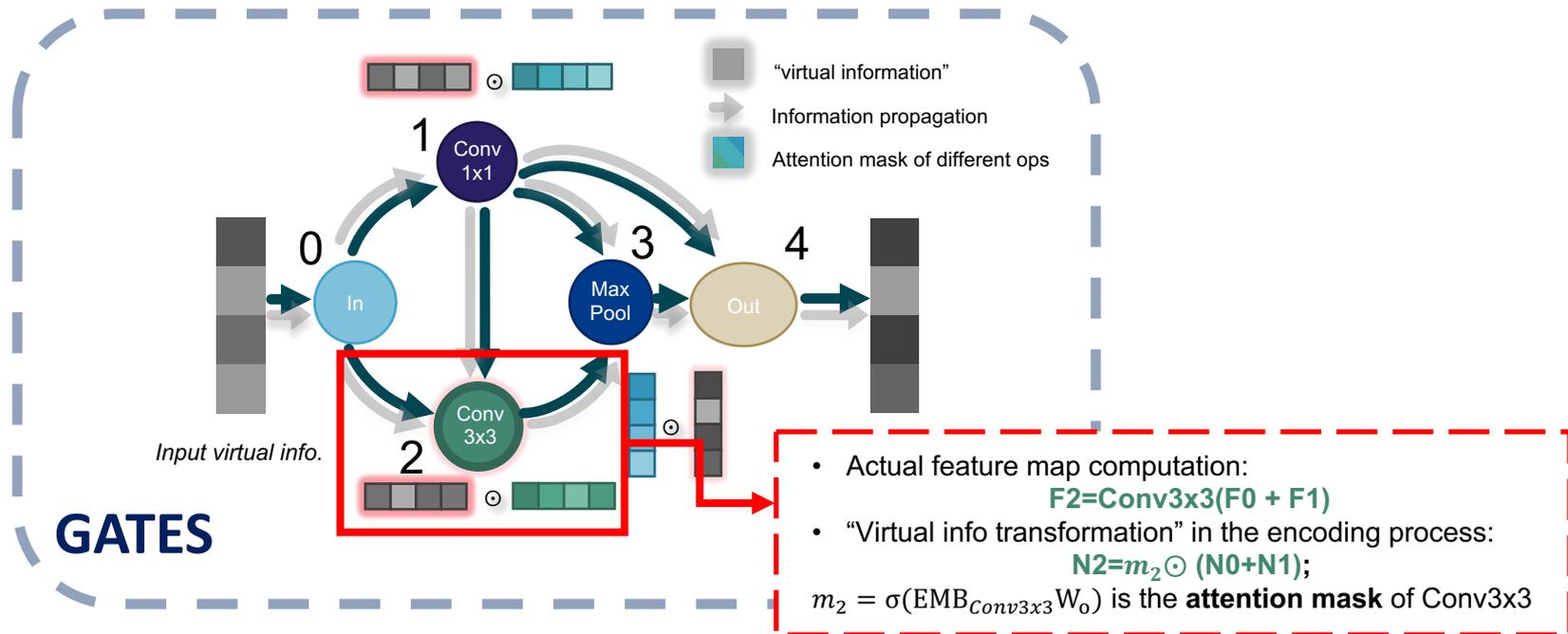


Should mind 2 aspects

- How to **encode**
- How to **train**

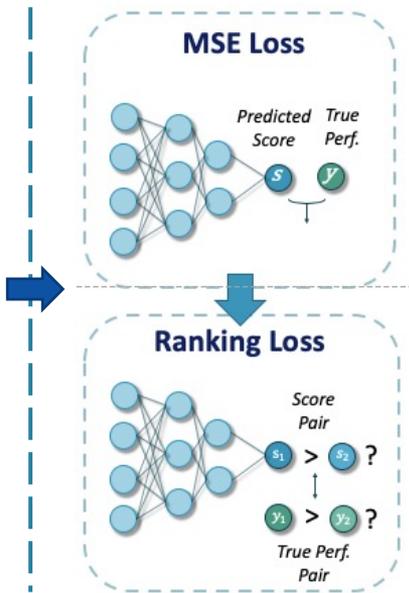
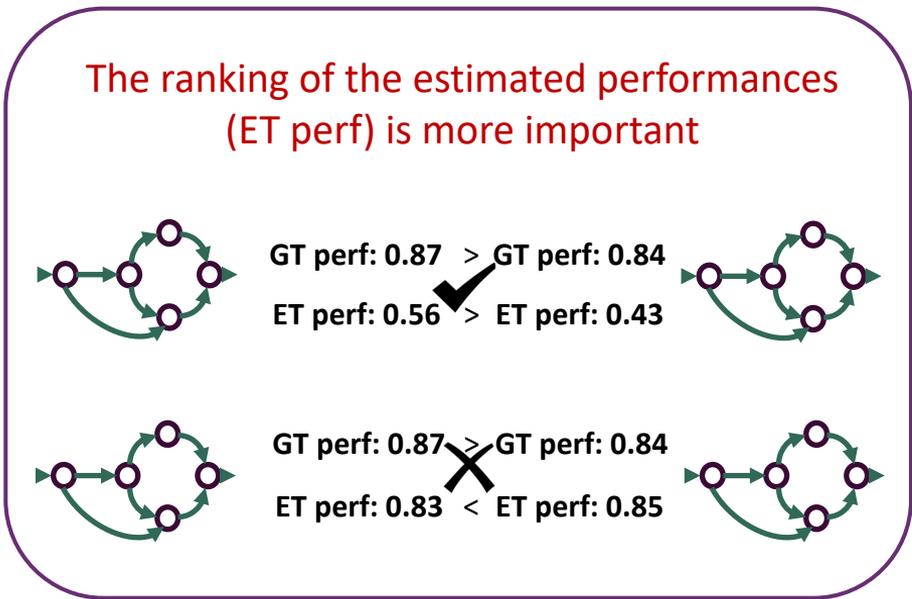
# GATES: How to Encode

- 🔑 **Mimic** the actual NN information flow
- 💡 Properly encode **computationally isomorphic** architectures to the same embedding



# GATES: How to Train

- Use **ranking loss** to train the predictor
- Ranking loss are better surrogate of **ranking performance** than regression loss



$$\sum_{j=1}^N (P(a_j) - y_j)^2$$

$$\sum_{j=1}^N \sum_{i, y_i > y_j} \max[0, m - (P(a_i) - P(a_j))]$$

# GATES: Results (on NAS-Bench-101)



- Ranking correlation (Kendall's Tau)

- Encoder comparison

Encoder	Proportions of 381262 training samples							
	0.05%	0.1%	0.5%	1%	5%	10%	50%	100%
MLP [21]	0.3971	0.5272	0.6463	0.7312	0.8592	0.8718	0.8893	0.8955
LSTM [21]	0.5509	0.5993	0.7112	0.7747	0.8440	0.8576	0.8859	0.8931
GCN (w.o. global node)	0.3992	0.4628	0.6963	0.8243	0.8626	0.8721	0.8910	0.8952
GCN (global node) [20]	0.5343	0.5790	0.7915	0.8277	0.8641	0.8747	0.8918	0.8950
<b>GATES</b>	<b>0.7634</b>	<b>0.7789</b>	<b>0.8434</b>	<b>0.8594</b>	<b>0.8841</b>	<b>0.8922</b>	<b>0.9001</b>	<b>0.9030</b>

**GATES outperform other encoders consistently, especially when there are few training samples**

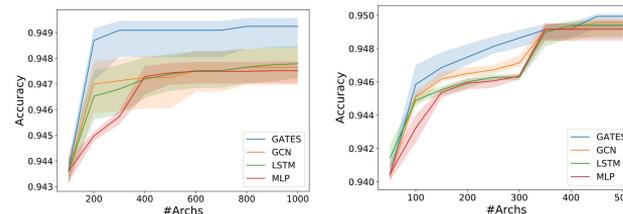
- Loss function comparison

Loss	Proportions of 381262 training samples							
	0.05%	0.1%	0.5%	1%	5%	10%	50%	100%
Regression (MSE) + GCN <sup>†</sup>	0.4536	0.5058	0.5587	0.5699	0.5846	0.5871	0.5901	0.5941
Regression (MSE) + GATES <sup>†</sup>	0.4935	0.5425	0.5739	0.6323	0.7439	0.7849	0.8247	0.8352
Pairwise (BCE)	0.7460	0.7696	0.8352	0.8550	0.8828	0.8913	0.9006	0.9042
Pairwise (Comparator)	0.7250	0.7622	0.8367	0.8540	0.8793	0.8891	0.8987	0.9011
Pairwise (Hinge)	0.7634	0.7789	0.8434	0.8594	0.8841	0.8922	0.9001	0.9030
Listwise (ListMLE)	0.7359	0.7604	0.8312	0.8558	0.8852	0.8897	0.9003	0.9009

**Ranking losses are better surrogate to ranking measures than regression losses**

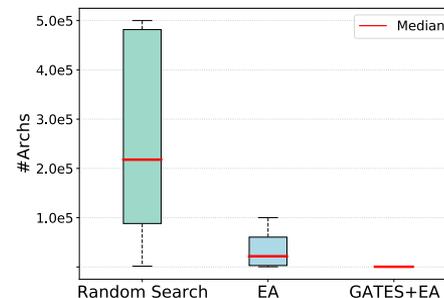
- Sample efficiency

- Encoder comparison



(a) RS inner search method ( $r = 500$ ) (b) EA inner search method ( $r = 100$ )

- Comparison with baseline search strategies



Median: 220k 24k 0.4k

**551.0x and 59.25x more efficient than RS/EA**

Ning et al., "A Generic Graph-based Neural Architecture Encoding Scheme for Predictor-based NAS", ECCV'20.

# Verification on HW-SW Co-design



🎯 Design **high accuracy and efficient** Process-In-Memory (PIM)-based system

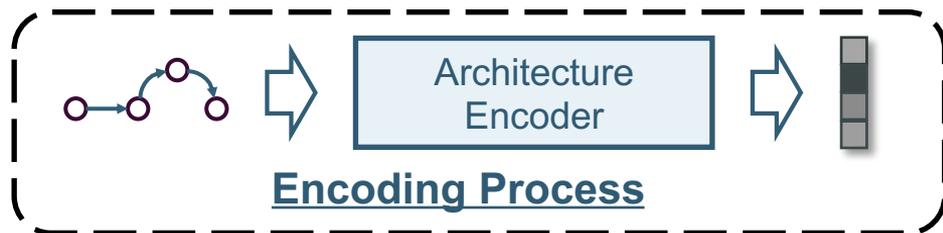
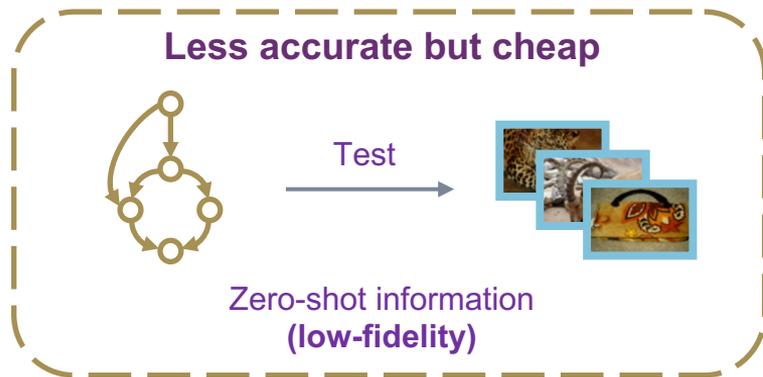
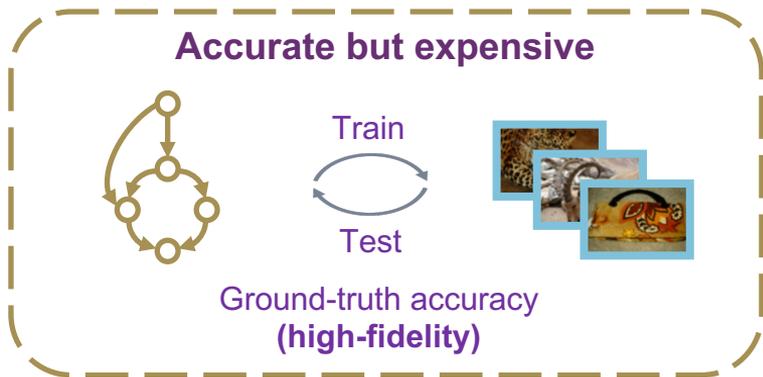
🔑 **Co-explore** the NN and PIM architectures based on **predictor-based NAS**

Method	NN accuracy	EDP ( $ms \times mJ$ )	Area ( $mm^2$ )	Search time ( $h$ )
NACIM [9]	73.9%	1.55	17.17	59
UAE [18]	83.0%	–	–	154
NAS4RRAM [10]	84.4%	–	–	289
CARS [16] (acc opt.)	88.0%	11.03	227.73	72
<i>Gibbon</i> (edp opt.)	84.6%	<b>0.24</b>	167.16	<b>7</b>
<i>Gibbon</i> (area opt.)	76.4%	1.00	<b>6.84</b>	<b>7</b>
<i>Gibbon</i> (acc opt.)	<b>88.3%</b>	14.33	186.32	<b>7</b>

- **0.2~10.7%** accuracy promotion
- **2.51×** area reduction
- **6.48×** EDP reduction
- **8.4~41.3×** search efficiency improvement

# DELE: How to Train

- Utilize **low-fidelity information** to help train the predictor
- More low-fidelity information **help learn better architecture representation**



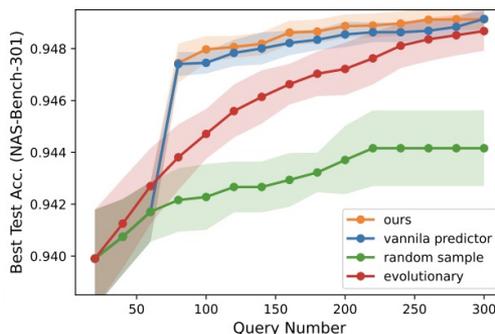
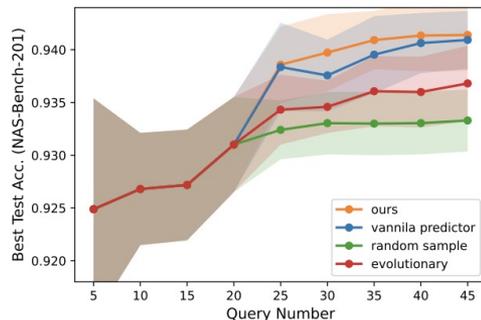
# DELE: Results

- Ranking correlation (Kendall's Tau)

Search Space	Encoder	Manner	Proportions of training samples				
			1%	5%	10%	50%	100%
NAS-Bench-201	GATES	Vanilla	0.7332 <sub>(0.0110)</sub>	0.8582 <sub>(0.0059)</sub>	0.8865 <sub>(0.0045)</sub>	0.9180 <sub>(0.0029)</sub>	0.9249 <sub>(0.0019)</sub>
		Ours	<b>0.8244</b> <sub>(0.0081)</sub>	<b>0.8948</b> <sub>(0.0021)</sub>	<b>0.9075</b> <sub>(0.0015)</sub>	<b>0.9216</b> <sub>(0.0019)</sub>	<b>0.9250</b> <sub>(0.0020)</sub>
NAS-Bench-201	LSTM	Vanilla	0.5692 <sub>(0.0087)</sub>	0.6410 <sub>(0.0018)</sub>	0.7258 <sub>(0.0053)</sub>	0.8765 <sub>(0.0010)</sub>	0.9000 <sub>(0.0008)</sub>
		Ours	<b>0.7835</b> <sub>(0.0062)</sub>	<b>0.8538</b> <sub>(0.0029)</sub>	<b>0.8683</b> <sub>(0.0015)</sub>	<b>0.8992</b> <sub>(0.0010)</sub>	<b>0.9084</b> <sub>(0.0010)</sub>
NAS-Bench-301	GATES	Vanilla	0.4160 <sub>(0.0450)</sub>	0.6752 <sub>(0.0088)</sub>	0.7354 <sub>(0.0044)</sub>	0.7693 <sub>(0.0041)</sub>	<b>0.7883</b> <sub>(0.0011)</sub>
		Ours	<b>0.5529</b> <sub>(0.0135)</sub>	<b>0.6830</b> <sub>(0.0038)</sub>	<b>0.7433</b> <sub>(0.0018)</sub>	<b>0.7752</b> <sub>(0.0026)</sub>	0.7842 <sub>(0.0022)</sub>
NAS-Bench-301	LSTM	Vanilla	0.4757 <sub>(0.0150)</sub>	0.6116 <sub>(0.0099)</sub>	0.6923 <sub>(0.0044)</sub>	0.7516 <sub>(0.0017)</sub>	0.7667 <sub>(0.0007)</sub>
		Ours	<b>0.4805</b> <sub>(0.0083)</sub>	<b>0.6405</b> <sub>(0.0035)</sub>	<b>0.7075</b> <sub>(0.0022)</sub>	<b>0.7544</b> <sub>(0.0028)</sub>	<b>0.7751</b> <sub>(0.0011)</sub>

DELE achieves better ranking performance across different search spaces, encoders and training proportions.

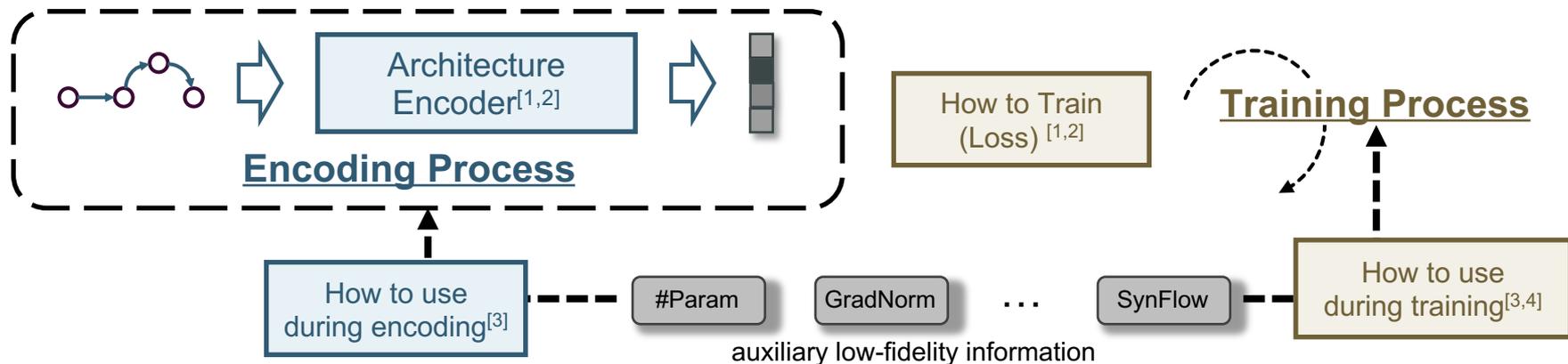
- Sample efficiency



DELE discovers better architectures with less query number.



# Summary of Architecture Encoding (Modeling)



How to train the encoder?

Depends on how we use the encoding!

How to utilize more information (what information)?

Incorporate the knowledge into encoding or through training.

## GATES<sup>[1]</sup>



Use encoding to rank architecture performance



Maximize ranking ability i.e., Minimize ranking loss

## CLOSE<sup>[2]</sup>



Use encoding to choose parameters for operations



Minimize parameter-sharing loss, train jointly with supernet

## DELE/GATES++<sup>[3,4]</sup>



Low-fidelity information can be beneficial for modeling



Utilize low-fidelity information in encoding and training

[1] Ning et al., "A Generic Graph-based Neural Architecture Encoding Scheme for Predictor-based NAS", ECCV/20.

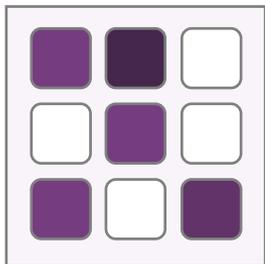
[2] Zhou\*, Ning\* et al., "CLOSE: Curriculum Learning On the Sharing Extent Towards Better One-shot NAS", ECCV/22.

[3] Ning et al., "A Generic Graph-based Neural Architecture Encoding Scheme with Multifaceted Information", TPAMI/23.

[4] Zhao\*, Ning\* et al., "Dynamic Ensemble of Low-fidelity Experts: Mitigating NAS Cold-Start", AAAI/23 Oral.

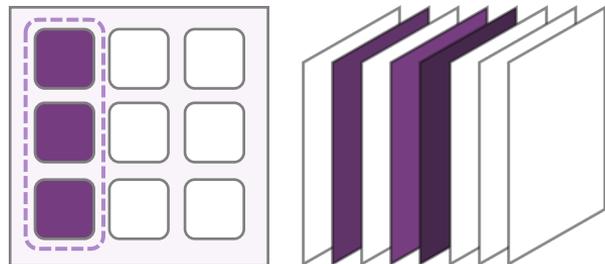
# Neural Network Pruning

- Definition: Given the **hardware-related constraints** (e.g., **latency, pruning ratio**), prune the **network structure** to satisfy the constraints and minimize the accuracy loss
- Pruning Granularity: Structured Pruning / Unstructured Pruning



**Unstructured Pruning**  
**Granularity: Weight**

(Example: Deep Compression [Han et. al. ICLR15])

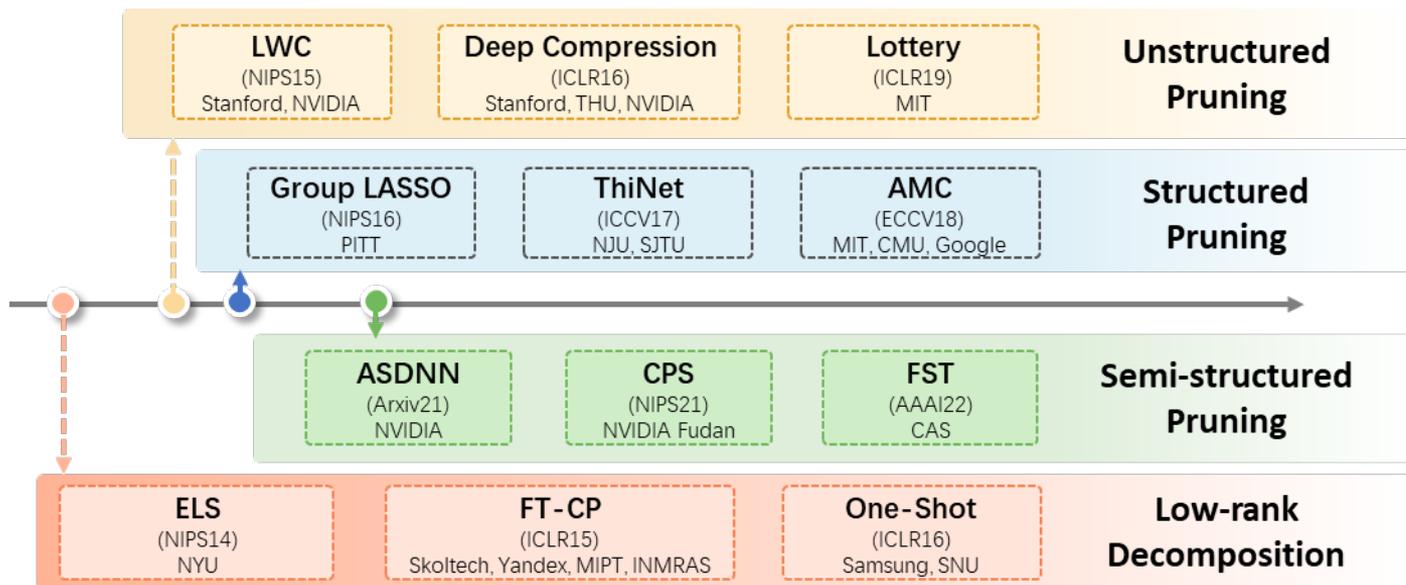


**Structured Pruning:**  
**Granularity: Channel/Group**  
(Example: AMC [Han et. al. ECCV18])

# Pruning Development Trends

- Trends

- Unstructured Pruning: Study on The Lottery Ticket Hypothesis
- Structured/Semi-structured Pruning: Hardware-Software Co-design for better trade-off between sparsity and performance

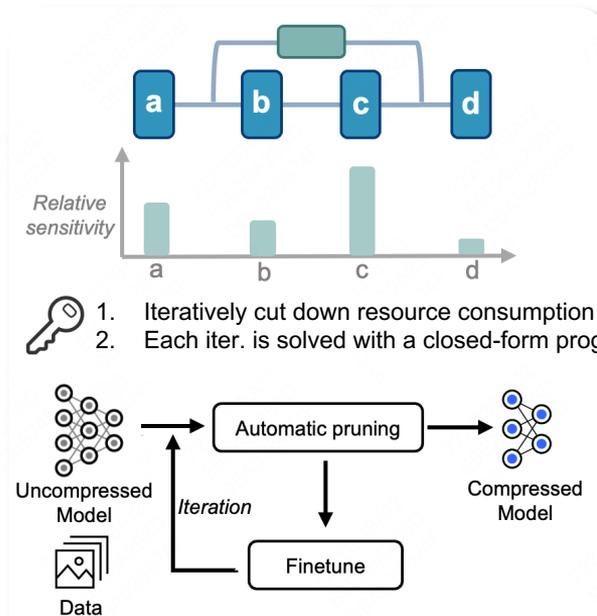


# Our Work on Effective Pruning

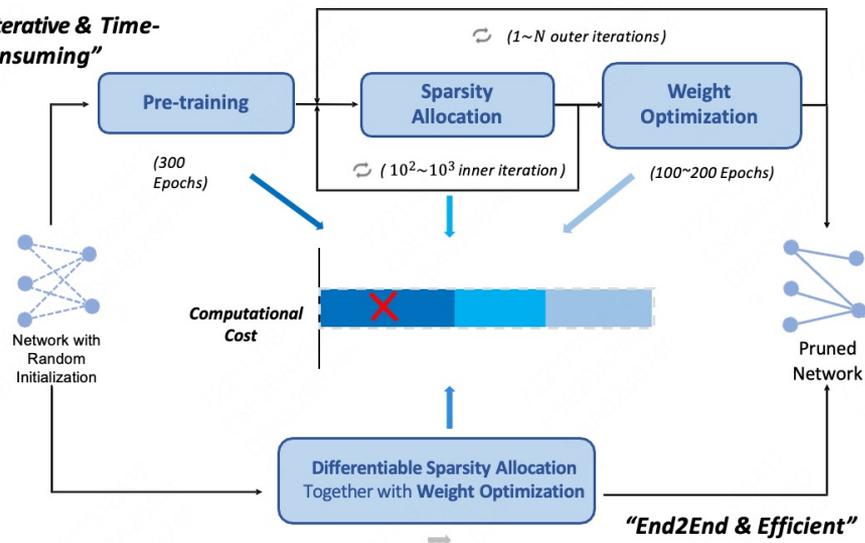
🎯 Design **controllable** pruning schemes

🔑 **Programming-based** iterative pruning<sup>[1,2]</sup>

🔑 **DSA: Differentiable Sparsity Allocation**<sup>[3]</sup>



*“Iterative & Time-consuming”*



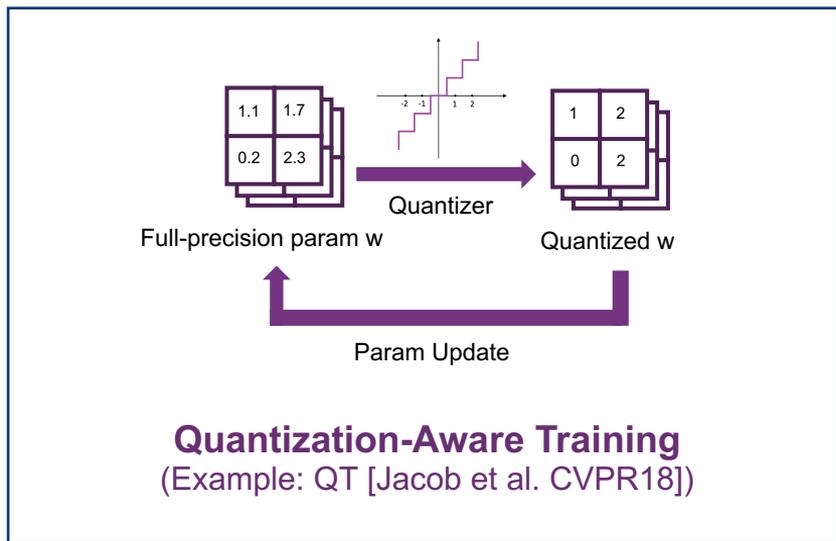
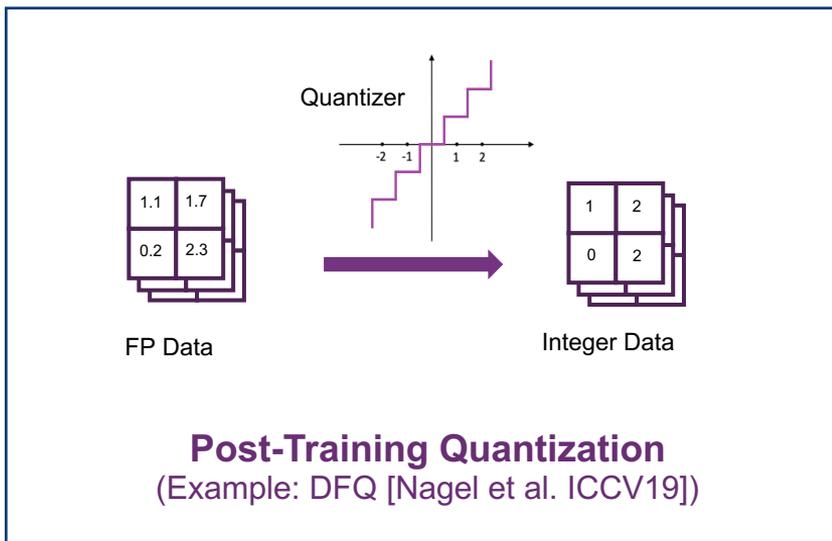
[1] Wang, **Ning** et al, Hardware Design and Software Practices for Efficient Neural Network Inference, *Low-Power Computer Vision*. Chapman and Hall/CRC 55-90. 2020.

[2] Shi\*, **Ning\***, Guo et al., Memory-Oriented Structural Pruning for Efficient Image Restoration, AAAI'23.

[3] **Ning** et al., DSA: More Efficient Budgeted Pruning via Differentiable Sparsity Allocation, ECCV'20.

# Neural Network Quantization

- Definition: Given the **hardware-related constraints** (e.g., **latency, bitwidth**), quantize FP32 tensors (**weights/activations**) to satisfy the constraints and minimize the accuracy loss
- Quantization Methods: Post-Training Quantization / Quantization-Aware Training

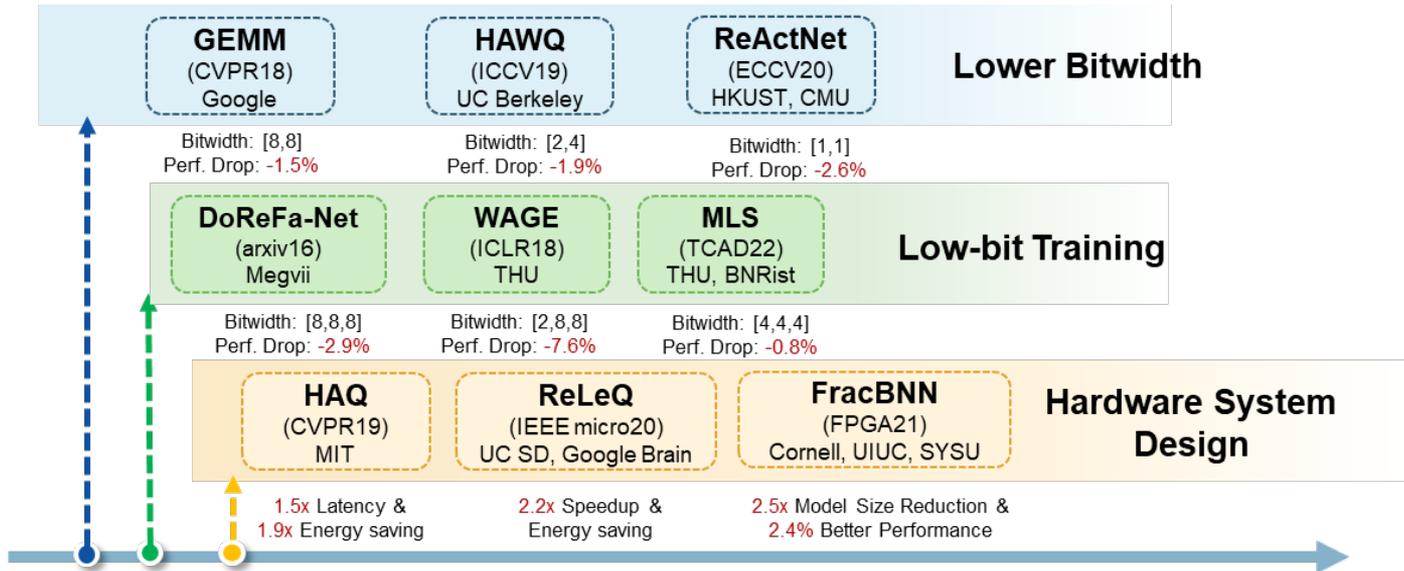




# Quantization Development Trends

## Trends

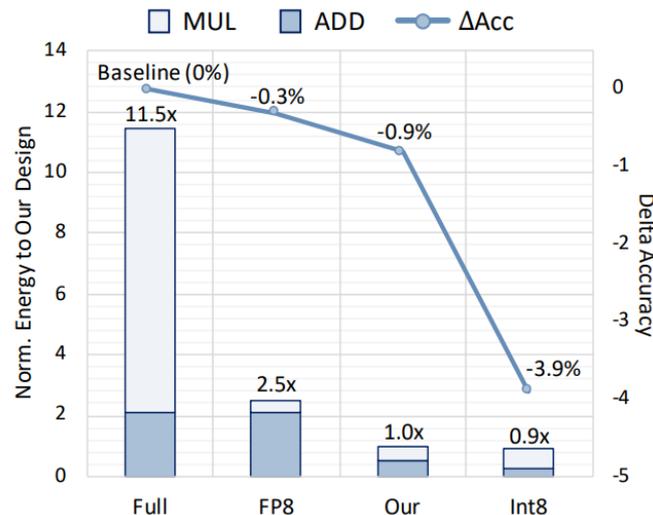
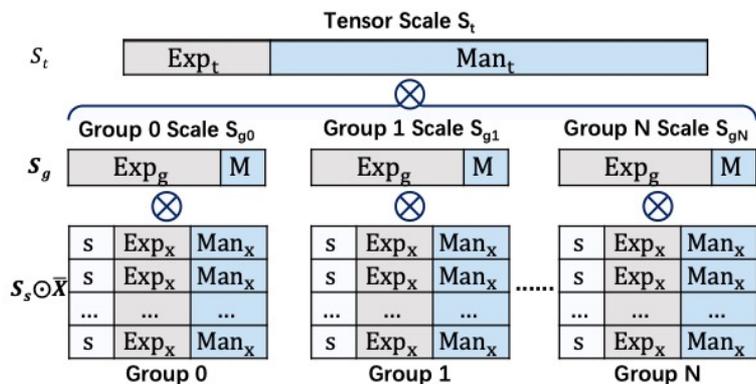
- Use **lower** bitwidth (8 → 4 → 2)
- Low-bit inference → **Training**
- **Hardware-Algorithm co-design** of the low-bit computing system



# Our Work on Effective Quantization

🎯 Achieve efficient **low-bit training**

🔑 Special format design that (1) is **hardware-friendly**; (2) maintains **representation range**



- 4-/2-bit training within **1%** accuracy loss
- **10.2x** energy efficiency than FP32

[1] Zhong, **Ning** et al., Exploring the Potential of Low-bit Training of Convolutional Neural Networks, TCAD 22.



## 目录 Contents

- 1 Why does efficient AI inference matters?
- 2 How to generally accelerate AI inference?
- 3 Research on specific applications
  - Low-level vision
  - Image generation
  - 3D perception



## 目录 Contents

- 1 Why does efficient AI inference matters?
- 2 How to generally accelerate AI inference?
- 3 Research on specific applications
  - Low-level vision with pruning
  - Image generation **with NAS**
  - 3D perception **with dynamic inference**



## 目录 Contents

- 1 Why does efficient AI inference matters?
- 2 How to generally accelerate AI inference?
- 3 Research on specific applications
  - Low-level vision
  - Image generation
  - 3D perception

# Background: Image Generation

- **Application field:** 2D image/3D assets generation task
  - Diverse Tasks & Efficiency Demands

“A shirt with the inscription: I love generative models”



Text-Image Generation



Image Editing



Image Super-resolution



3D Generation

## Diverse Tasks

**Low Memory:** Deploying models on PC or Mobile Phone (**500M~1G runtime memory**) becomes popular in many application scenario.



**Firm Real-time:** The generation speed of the AI model should reach **1s/img**

- Stable Diffusion has been used as a plug-in to help paint in Photoshop<sup>[1]</sup>, which demands fast generation speed to satisfy the users' frequent requests.

## Efficiency Demands

[1] Photoshop Stable Diffusion Plugin, <https://github.com/AbdullahAlfaraj/Auto-Photoshop-StableDiffusion-Plugin>

# Application Challenges



## Challenges of Diffusion Models

### Application Demands

- **Low Memory Cost: 500M~1G** runtime memory cost on edge devices
- **Low Latency: Achieve 1s/img (firm real-time)** speed on edge devices

### Future Trends

- **Larger-scale:** Larger resolution, more training data
- **More tasks & Multi-Modal :** More conditioning types and more data modal.

### Efficiency of current methods



Nvidia A100

Sampler*	Memory (512x512)	Latency (512x512)
DDIM <sup>[2]</sup>	~12G	7~14s/img (100-200 NFE)
PNDM <sup>[3]</sup>		3.5~7s/img (50-100 NFE)
DPM-Solver <sup>[4]</sup>		1.4~3.5s/img (20-50 NFE)

\* Stable diffusion inference using different samplers.



**Call for 10x~100x Efficiency Improvement**

[1] Rombach et al., High-Resolution Image Synthesis with Latent Diffusion Models, CVPR'22.

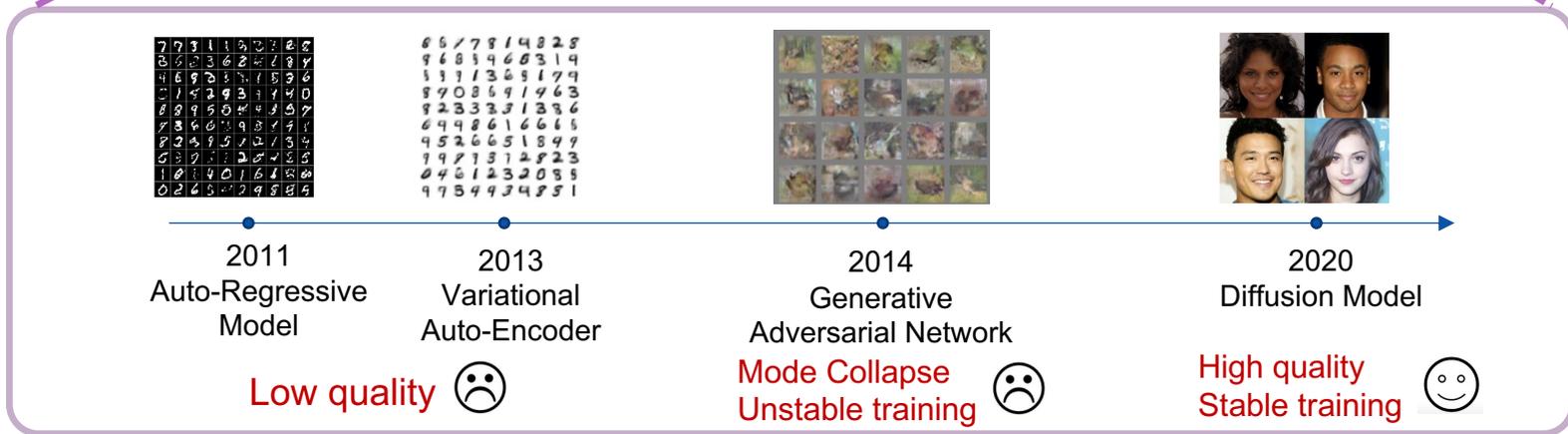
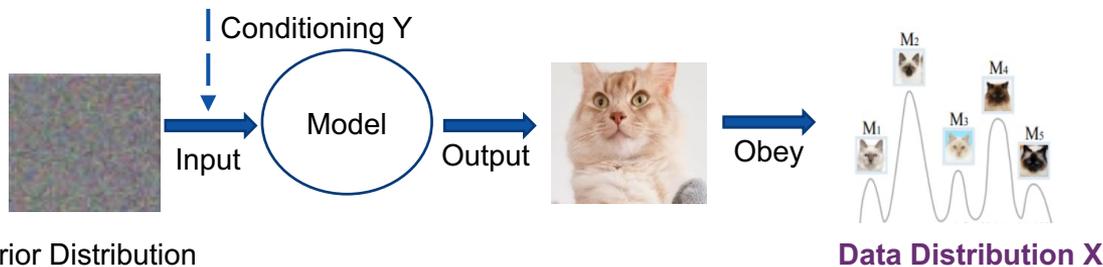
[2] Song et al., Denoising Diffusion Implicit Models, ICLR'21.

[3] Liu et al., Pseudo Numerical Methods for Diffusion Models on Manifolds, ICLR'22.

[4] Lu et al., DPM-Solver: A Fast ODE Solver for Diffusion Probabilistic Model Sampling in Around 10 Steps, NeurIPS'22.

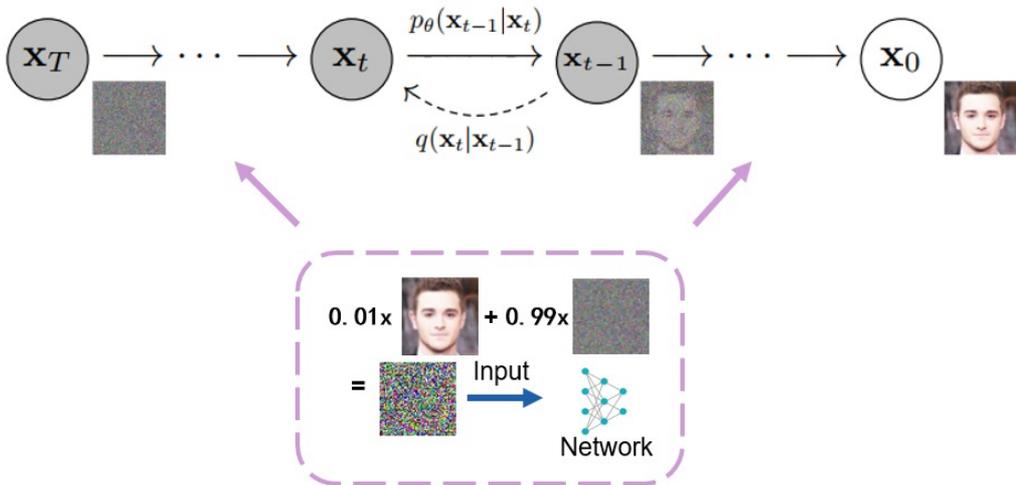
# Background: Generative Model

- Generative model aims at learning the **data distribution**



# Background & Motivation

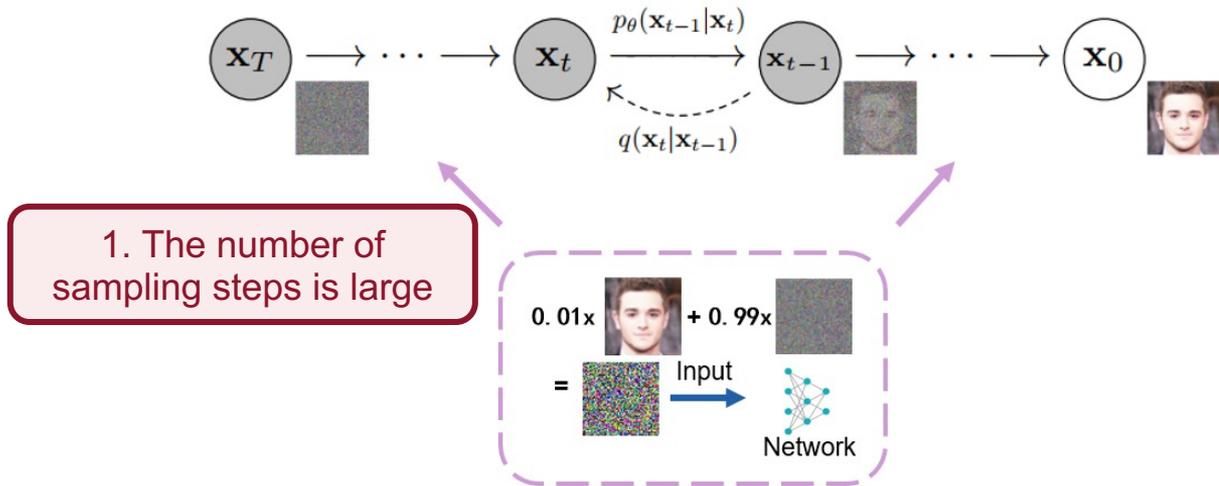
- Diffusion in a nutshell: learn to **iteratively denoise a gaussian noise** for generation



Sampler	NFE	Inference time per step (NVIDIA A100, BS=1)	Total Time (NVIDIA A100, BS=1)
DDIM <sup>[2]</sup>	100~200	~0.07s	7-14s/image
PNDM <sup>[3]</sup>	50~100		3.5~7s/image
DPM-Solver <sup>[4]</sup>	20~50		1.4~3.5s/image

# Background & Motivation

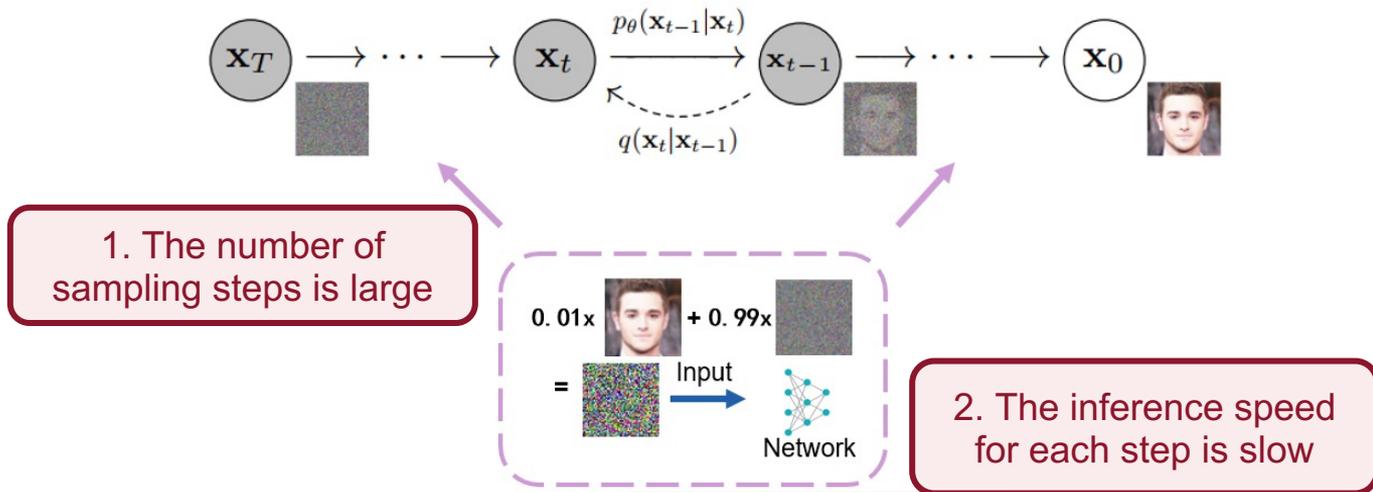
- Diffusion in a nutshell: learn to **iteratively denoise a gaussian noise** for generation



Sampler	NFE	Inference time per step (NVIDIA A100, BS=1)	Total Time (NVIDIA A100, BS=1)
DDIM <sup>[2]</sup>	100~200	~0.07s	7-14s/image
PNDM <sup>[3]</sup>	50~100		3.5~7s/image
DPM-Solver <sup>[4]</sup>	20~50		1.4~3.5s/image

# Background & Motivation

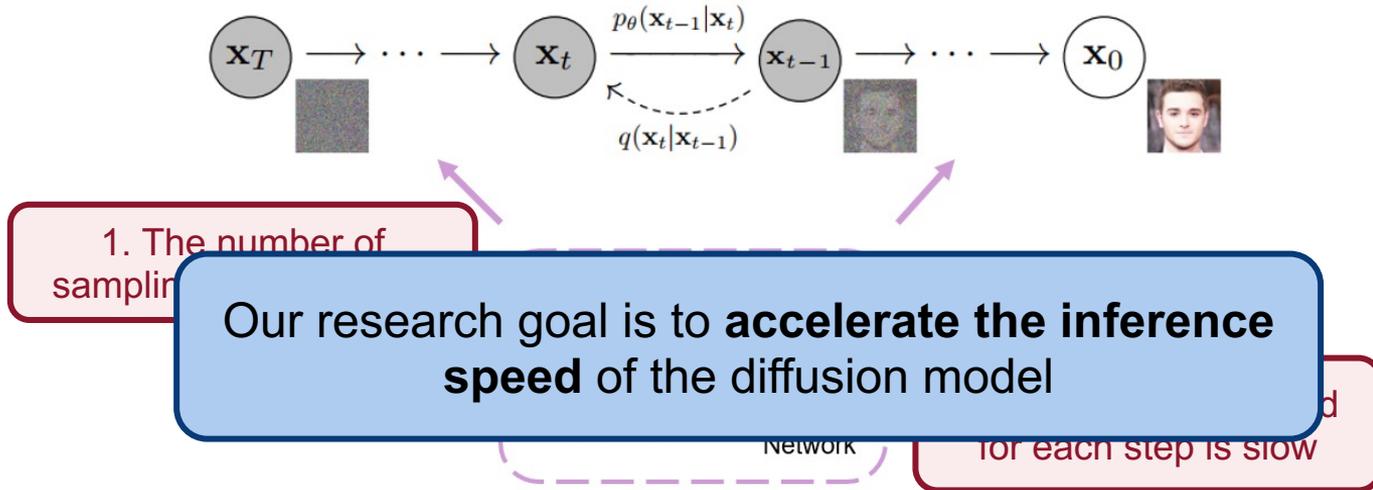
- Diffusion in a nutshell: learn to **iteratively denoise a gaussian noise** for generation



Sampler	NFE	Inference time per step (NVIDIA A100, BS=1)	Total Time (NVIDIA A100, BS=1)
DDIM <sup>[2]</sup>	100~200	~0.07s	7-14s/image
PNDM <sup>[3]</sup>	50~100		3.5~7s/image
DPM-Solver <sup>[4]</sup>	20~50		1.4~3.5s/image

# Background & Motivation

- Diffusion in a nutshell: learn to **iteratively denoise a gaussian noise** for generation



Sampler	NFE	Inference time per step (NVIDIA A100, BS=1)	Total Time (NVIDIA A100, BS=1)
DDIM <sup>[2]</sup>	100~200	~0.07s	7-14s/image
PNDM <sup>[3]</sup>	50~100		3.5~7s/image
DPM-Solver <sup>[4]</sup>	20~50		1.4~3.5s/image

# Idea and Oracle Experiments

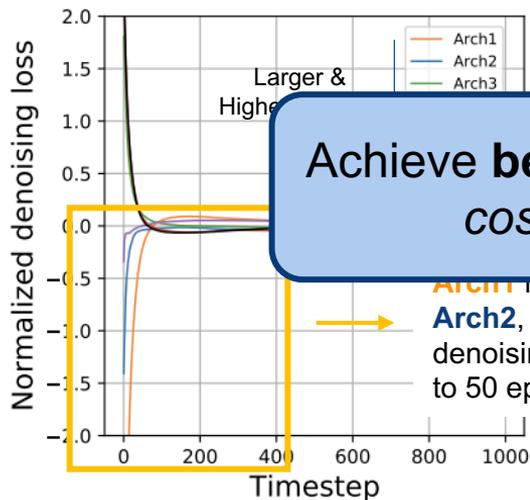
- Existing SOTA methods adopt the **same model** during **all the diffusion steps**
  - Use different models in different steps to reduce the computation demand?**

Observation1

**Small models outperform** large models at **some steps** in terms of denoising loss

Observation2

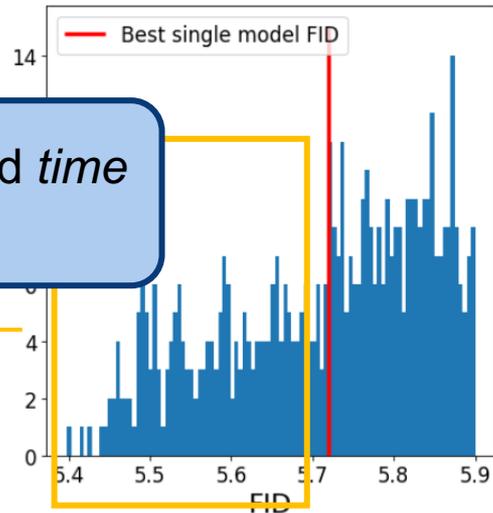
Some **mixtures of large and small models** get **better** performance than **only use large model**



Achieve **better trade-off** between *quality* and *time cost* by optimizing **model schedule**

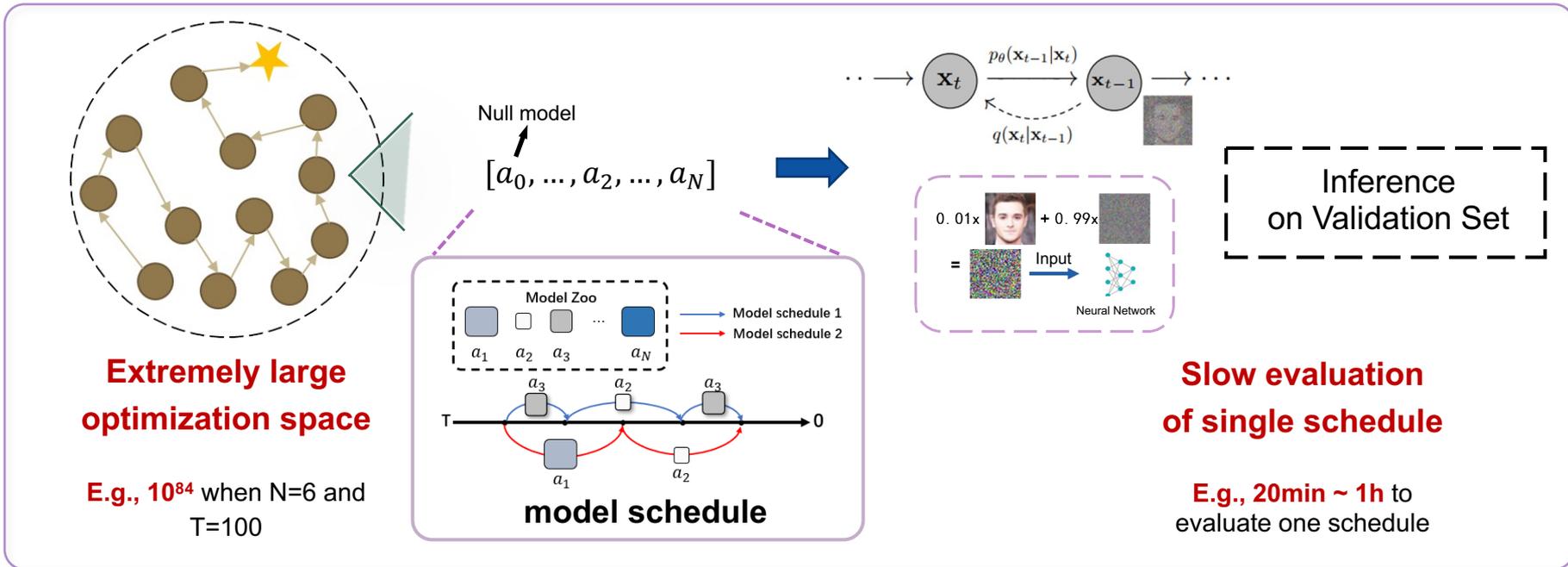
Arch1 is smaller than Arch2, but obtains smaller denoising loss between 1 to 50 epochs.

Mixing small and large models can possibly get better FID than only use large model



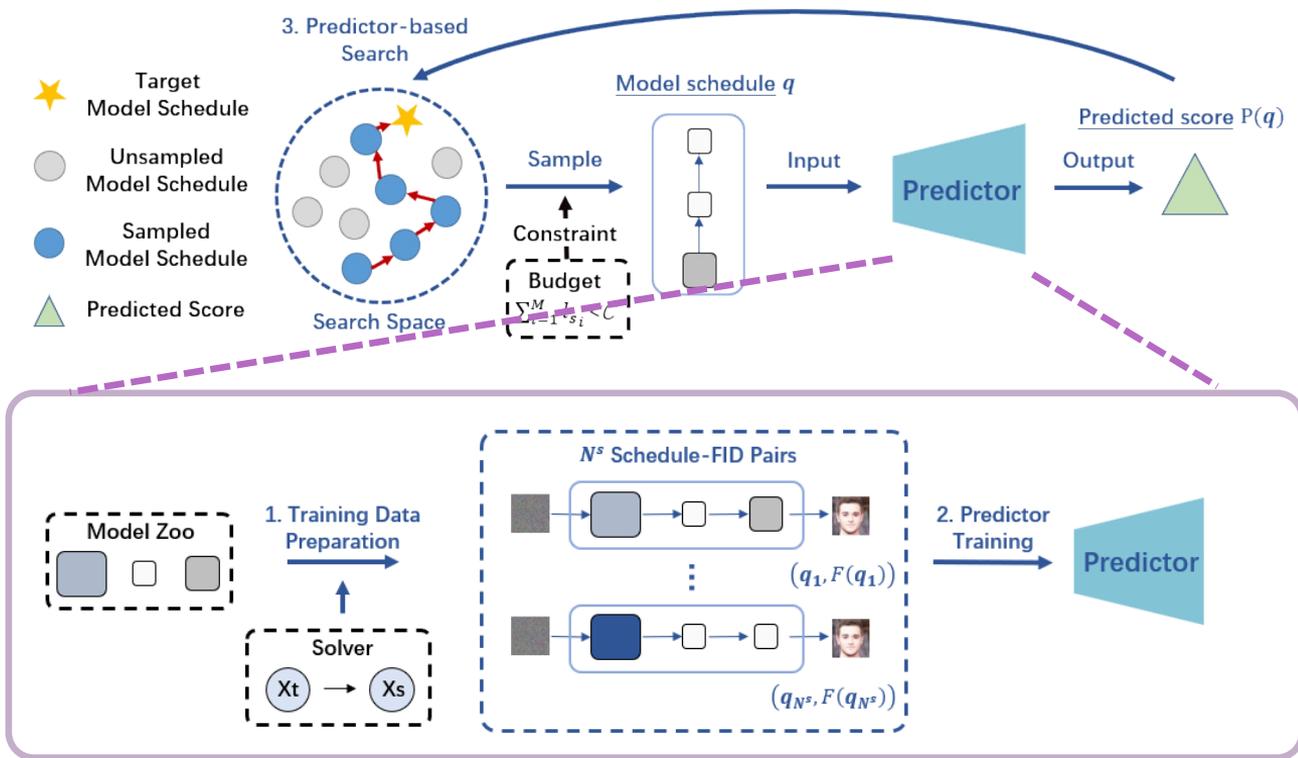
# Problem Formulation and Challenges

- Formulation: Given a **pre-trained model zoo** and the **generation time constraint**, search for the **model schedule** to satisfy the constraint and maximize the quality score (i.e., FID)
- Challenges: **Slow search process**



# Method

- **Predictor-based search** for a good model schedule satisfying the constraint

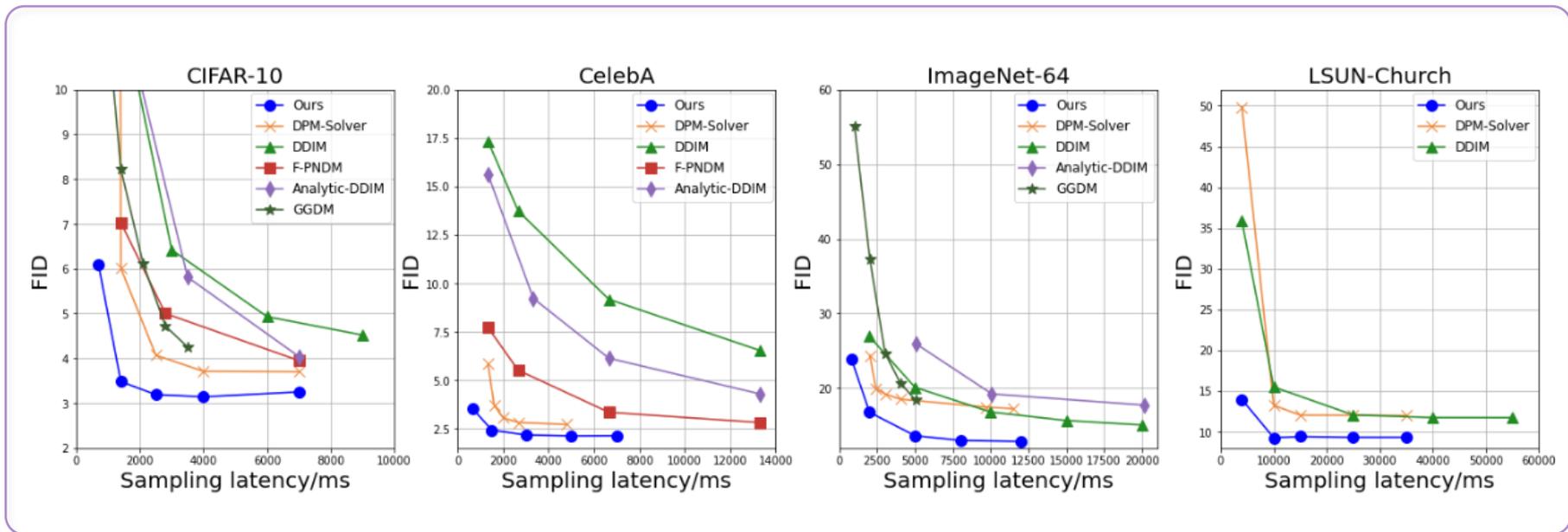


# Results



## • Experimental Results

- Accelerate generation process by **2-5x** (on A100) without sacrificing the performance



## • Experimental Results

- Accelerate generation process by **2-5x** (on A100) without sacrificing the performance
- Achieve **better FID** than *Stable Diffusion* by optimizing the schedule

Budget/NFE	Baseline Type		Ours
	(1)	(2)	
9	13.01	13.01	<b>12.90</b>
12	12.11	11.37	<b>11.34</b>
15	11.92	11.13	<b>10.72</b>
18	11.88	11.13	<b>10.68</b>
24	11.81	11.13	<b>10.57</b>

Table 3. FID on MS-COCO 256×256. All FIDs in the table are calculated between 30k images in validation set and 30k generated images guided with the same captions.

# Results



- **Demonstration** of generated images with the same generation time

**Previous  
SOTA Method**



**LSUN-Church**  
(Budget: 4s, bs=64)

**Our Method**



**Stable Diffusion**  
(Budget: 15 NFE)



## Unconditional Generation

- Common
  - Low budgets: Small models & **Sufficient steps**
  - High budgets: **Large models**
- Varies with dataset
  - CIFAR-10/CelebA/ImageNet-64 (32x32, 64x64):
    - **Larger** models & **Denser** steps for lower t
  - LSUN-Church (256x256):
    - **Smaller** models & **Sparser steps** for lower t

## Conditional Generation with Stable Diffusion (512x512)

- Mixing **checkpoints with the same complexity but different functionalities** is still beneficial for the efficiency-quality trade-off
- **1-st solver** is preferred for lower t

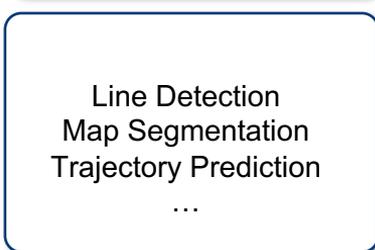
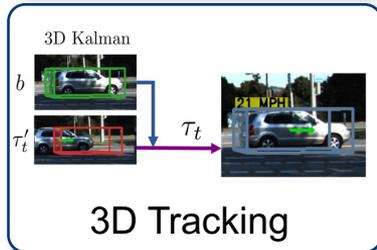
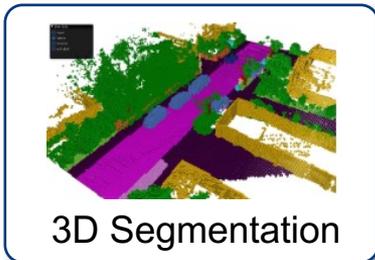
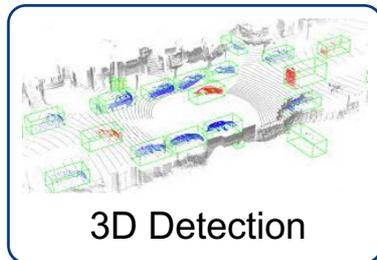


## 目录 Contents

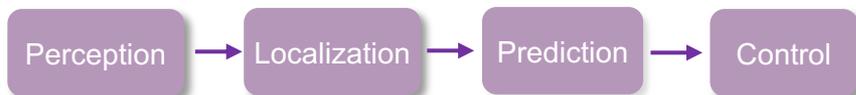
- 1 Why does efficient AI inference matters?
- 2 How to generally accelerate AI inference?
- 3 Research on specific applications
  - Low-level vision
  - Image generation
  - 3D perception

# Background: 3D Perception

- **Application field:** 3D scene understanding for autonomous driving
  - Diverse Tasks & Efficiency Demands



## Diverse Tasks

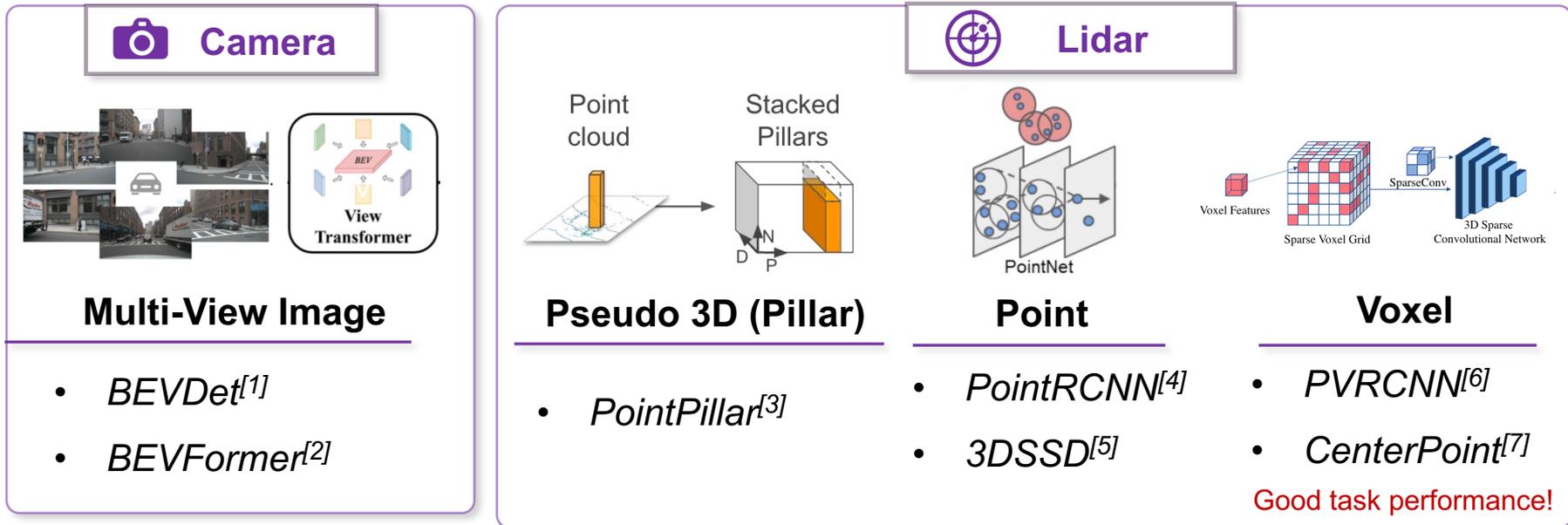


**Real-time:** According to different scenarios (e.g., APA, High-way, City) and functionalities, the **perception** system should reach **10~30Hz** (e.g., Tesla vision can reach 36Hz)

## Efficiency Demands

# Background: 3D Perception

- **Application field:** 3D scene understanding for autonomous driving
  - Different schemes



[1] Huang et al., : BEVDet: High-performance Multi-camera 3D Object Detection in Bird-Eye-View, Arxiv.

[2] Li et al., BEVFormer: Learning Bird's-Eye-View Representation from Multi-Camera Images via Spatiotemporal Transformers, ECCV22.

[3] Alex et al., PointPillars: Fast Encoders for Object Detection from Point Clouds, CVPR19.

[4] Shi et al., PointRCNN: 3D Object Proposal Generation and Detection from Point Cloud, CVPR19.

[5] Yang et al., 3DSSD: Point-based 3D Single Stage Object Detector, CVPR20.

[6] Shi et al., PV-RCNN: Point-Voxel Feature Set Abstraction for 3D Object Detection, CVPR20.

[7] Yin et al., Center-based 3D Object Detection and Tracking, CVPR21.

# Application Challenges

- Application challenges of 3D Perception Methods

## Application Demands

- **Low Latency:** Perception **10~30FPS**
- **Low Power:** Resource-constrained

## Efficiency of current methods



Nvidia Xavier



Nvidia RTX2080

## Futu

Our research goal is to **accelerate the inference speed** of the current 3D perception methods

- **Larger-scale:** Larger scene, more data
- **Time Dimension:** Process multiple frames in point cloud video
- **Multi-Modal:** Process multiple modality

	FPS
(Pillar-based)	~60
[1] (Point-based)	~8
PVRCNN <sup>[6]</sup> (Voxel-based)	~8
BEVDet <sup>[1]</sup> (Image-based)	~10



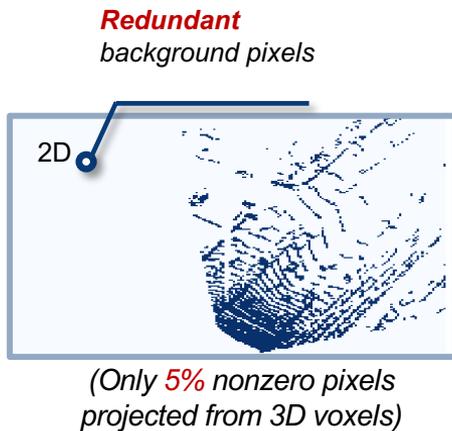
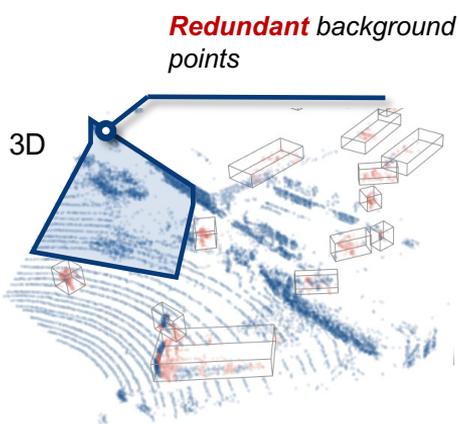
**Call for ~10x Efficiency Improvement**

# Idea

- Exploit the **Spatial Data Redundancy** with **Adaptive Inference**

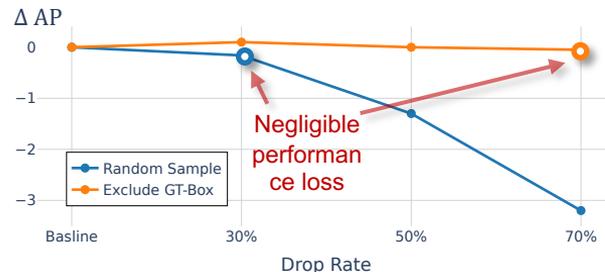
## Idea

3D detector needs to process **large scale** point cloud  
**Spatial redundancy** exists for both the 3D and BEV space



## Oracle

Validate **the spatial redundancy exists**,  
 and correctly dropping redundant part  
 brings **negligible performance drop**.



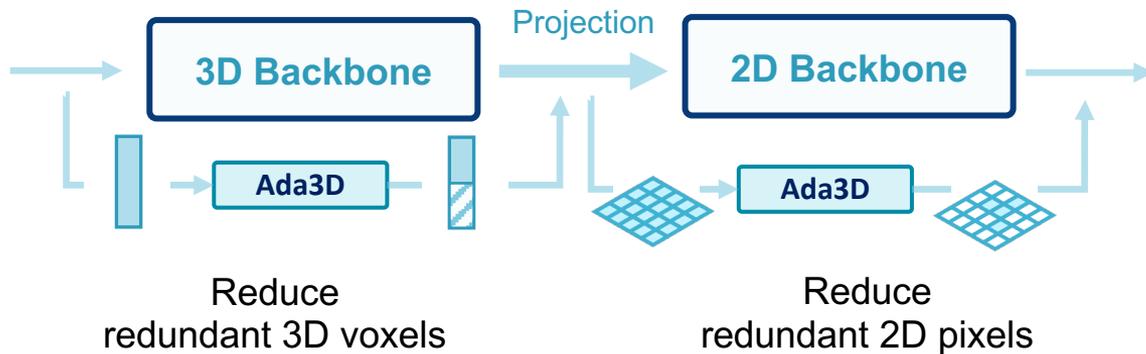
(\* Upper bound for adaptive inference)

# Idea

- Exploit the **Spatial Data Redundancy** with **Adaptive Inference**

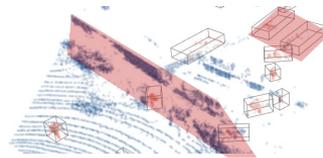
## Idea

- Reduce spatial redundancy by **adaptively skipping** computation and saving for redundant 3D/2D features during inference.
- A **share lightweight predictor** could effectively identify redundant features for each layer.

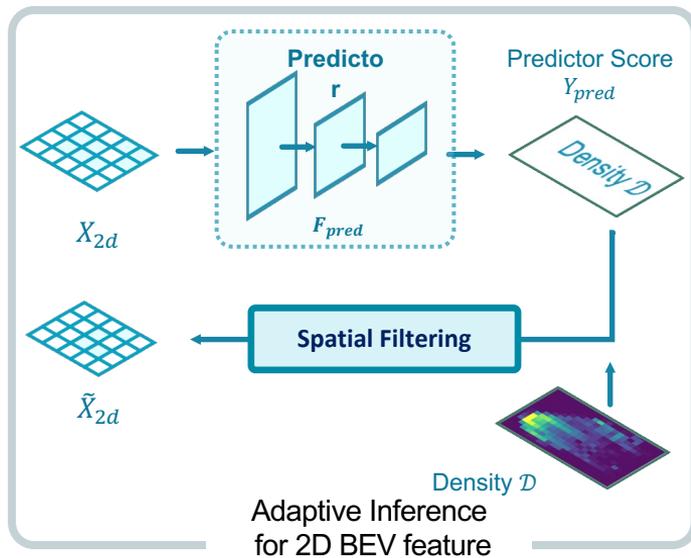
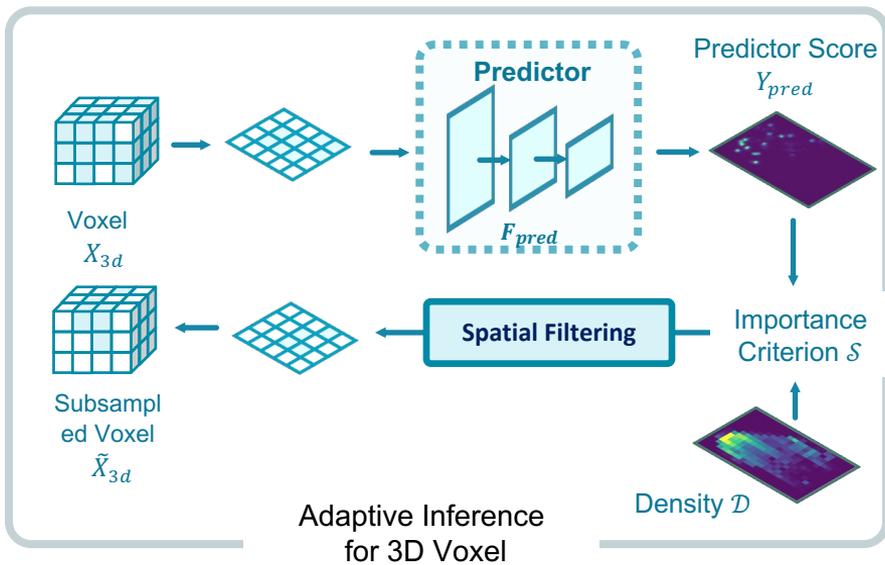


## • Adaptive inference

- Feature-based Importance Predictor
- Density-guided Spatial Filtering



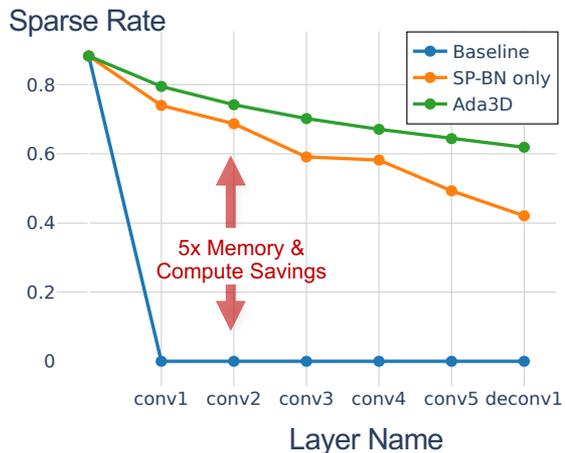
Local Point cloud have **higher density**



# Methods

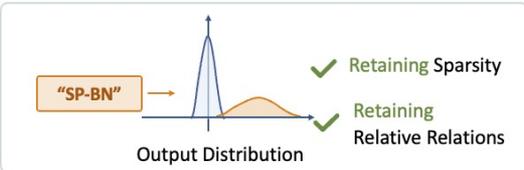
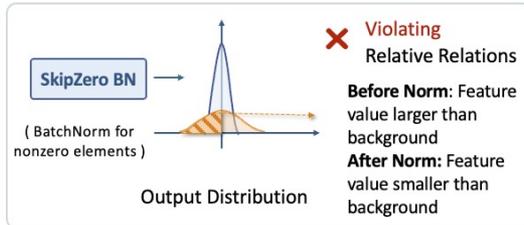
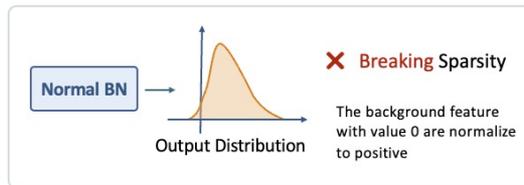
## • Sparsity-Preserving Batch Normalization

- Apply batch normalization for nonzero elements without mean subtraction



(Only 5% nonzero pixels projected from 3D voxels)

Most existing literatures treat the 2D BEV space as **dense** features



**SP-BN retains the sparsity of 2D BEV features**

# Experimental Results



## Experiments on KITTI Test

- Ada3D achieves **comparable performance** with baseline methods with **5~10x FLOPs and memory** compress rate.

Table 1: Performance comparison of Ada3D and other methods on KITTI *test set*. The “Ada3D-B” and “Ada3D-C” are centerpoint models optimized by Ada3D with different drop rates.

Mehod	FLOPs Opt.	Mem Opt.	mAP (Mod.)	3D Car (IoU=0.7)			3D Ped. (IoU=0.5)			3D Cyc. (IoU=0.5)		
				Easy	Mod.	Hard	Easy	Mod.	Hard	Easy	Mod.	Hard
VoxelNet [30]	-	-	49.05	77.47	65.11	57.73	39.48	33.69	31.50	61.22	48.36	44.37
SECOND [21]	-	-	57.43	84.65	75.96	68.71	45.31	35.52	33.14	75.83	60.82	53.67
PointPillars [8]	-	-	58.29	82.58	74.31	68.99	51.45	41.92	38.89	77.10	58.65	51.92
SA-SSD [4]	-	-	-	88.75	79.79	74.16	-	-	-	-	-	-
TANet [12]	-	-	59.90	84.39	75.94	68.82	53.72	44.34	40.49	75.70	59.44	52.53
Part-A <sup>2</sup> [17]	-	-	61.78	87.81	78.49	73.51	53.10	43.35	40.06	79.17	63.52	56.93
SPVCNN [20]	-	-	61.16	87.80	78.40	74.80	49.20	41.40	38.40	80.10	63.70	56.20
PointRCNN [16]	-	-	57.95	86.96	75.64	70.70	47.98	39.37	36.01	74.96	58.82	52.53
3DSSD [24]	-	-	55.11	87.73	78.58	72.01	35.03	27.76	26.08	66.69	59.00	55.62
IA-SSD [28]	-	-	60.30	88.34	80.13	75.10	46.51	39.03	35.60	78.35	61.94	55.70
CenterPoint [25]	-	-	59.96	88.21	79.80	76.51	46.83	38.97	36.78	76.32	61.11	53.62
CenterPoint-Pillar [25]	-	-	57.39	84.76	77.09	72.47	44.07	37.80	35.23	75.17	57.29	50.87
CenterPoint (Ada3D-B)	<b>5.26×</b>	<b>4.93×</b>	<b>59.85</b>	87.46	79.41	75.63	46.91	39.11	36.43	76.09	61.04	53.73
CenterPoint (Ada3D-C)	<b>9.83×</b>	<b>8.49×</b>	<b>57.72</b>	82.52	74.98	69.11	43.66	38.23	34.80	75.27	59.96	52.14

# Experimental Results

## • Experiments on KITTI Val

- Ada3D-A: **Improve the performance** with 80% 2D reduction
- Ada3D-B: **Without performance loss**, reduce 40% 3D features, 80% 2D features, compress FLOPs and memory cost by 5x
- Ada3D-C: **With moderate performance loss**, reduce 60% 3D features. 90% 2D features, improve model FLOPs and memory by **an order of magnitude**

Table 3: **Ablation studies and quantitive efficiency improvements of different Ada3D models on KITTI val.** “IP” stands for “importance predictor”, “DG” for “density-guided spatial filtering”, “SP-BN” for “sparsity preserving batch normalization”. The “FLOPs” and “Mem.” calculates the normalized resource consumption of the optimized model.

Method	Technique			FLOPs		Mem.		mAP (Mod.)	Car Mod. (IoU=0.7)	Ped. Mod. (IoU=0.5)	Cyc. Mod. (IoU=0.5)
	IP	DG	SP-BN	3D	2D	3D	2D				
CenterPoint	-	-	-	1.00	1.00	1.00	1.00	66.1	79.4 (-)	53.4 (-)	65.5 (-)
CenterPoint (SP-BN)	-	-	✓	1.00	0.49	1.00	0.45	66.0	79.1 (-0.3)	53.3 (-0.1)	65.6 (+0.1)
CenterPoint (Ada3D-A)	✓	✓	✓	1.00	0.22	1.00	0.25	66.4	79.5 (+0.1)	53.6 (+0.2)	66.1 (+0.6)
CenterPoint (Ada3D-B)	✓	✓	✓	0.66	0.18	0.68	0.17	66.1	79.1 (-0.3)	54.0 (+0.6)	65.3 (-0.3)
CenterPoint (Ada3D-B w.o. DG)	✓	-	✓	0.64	0.18	0.66	0.16	65.1	78.8 (-0.6)	51.6 (-1.8)	64.9 (-0.6)
CenterPoint (Ada3D-C)	✓	✓	✓	0.39	0.08	0.43	0.07	65.4	77.6 (-1.8)	53.5 (+0.2)	65.1 (-0.4)

# Experimental Results

- **Experiments on nuScenes and ONCE**
  - With less than **1% mAP loss**, compress the FLOPs and memory cost by **2~4x**
  - Achieve higher compression rate with less performance loss than **model-level compression** methods (e.g. SPSS-pruning & Channel Scaling)

Table 1: **Performance comparison of Ada3D on ONCE val set.** The results are taken from the OpenPCDet [5] implementation.

Method	<i>FLOPs</i> <i>Opt.</i>	<i>Mem.</i> <i>Opt.</i>	mAP	Veh.	Ped.	Cyc
PointRCNN [3]	-	-	28.74	52.09	4.28	29.84
PointPillar [1]	-	-	44.34	68.57	17.63	46.81
SECOND [6]	-	-	51.89	71.16	26.44	58.04
PVRCNN [2]	-	-	53.55	77.77	23.50	59.37
CenterPoint [7]	-	-	63.99	75.69	49.80	66.48
CenterPoint (with Ada3D)	2.32×	2.61×	62.68	73.43	49.09	65.53

Table 2: **Performance comparison of Ada3D on Nuscenes val set.** The “SPSS-Conv” model applies pruning for the 3D sparse convolution only, and the “CenterPoint-Mini” uses the 2D backbone with half the usual width.

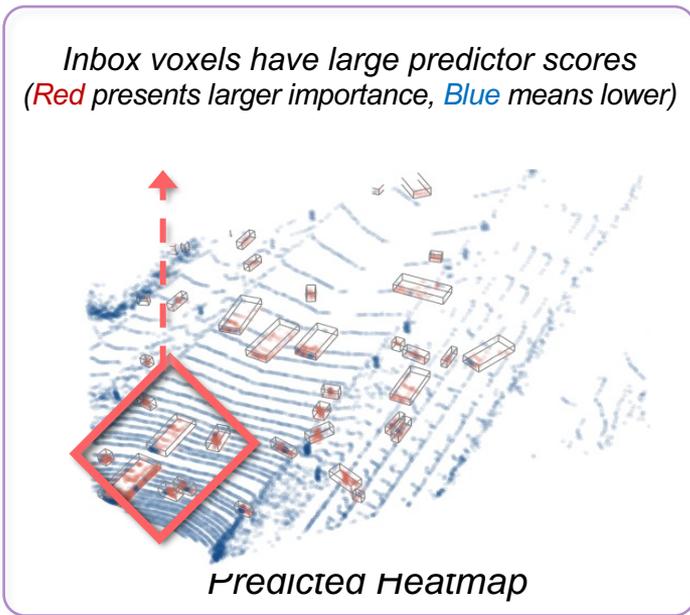
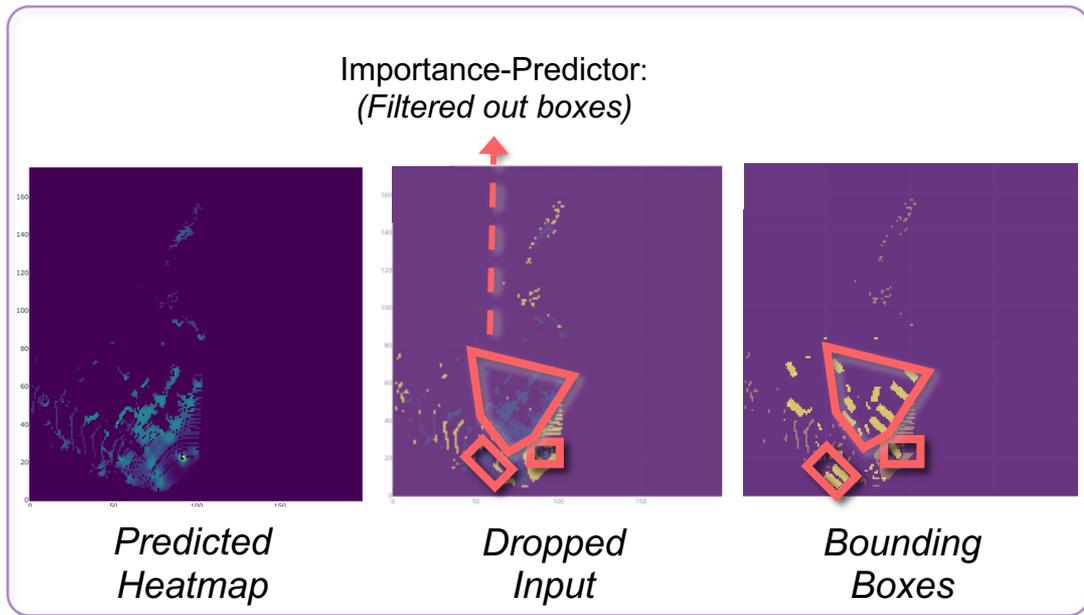
Method	<i>FLOPs</i> <i>Opt.</i>	<i>Mem.</i> <i>Opt.</i>	mAP	NDS
PointPillar [8]	-	-	44.63	58.23
SECOND [21]	-	-	50.59	62.29
CenterPoint-Pillar [25]	-	-	50.03	60.70
CenterPoint [25] ( <i>voxel=0.1</i> )	-	-	55.43	64.63
CenterPoint-Ada3D ( <i>voxel=0.1</i> )	2.32×	2.61×	54.80	63.53
CenterPoint [25] ( <i>voxel=0.075</i> )	-	-	59.22	66.48
SPSS-Conv [11] ( <i>voxel=0.075</i> )	1.14×	1.14×	57.80	65.69
CenterPoint-Mini [25] ( <i>voxel=0.075</i> )	2.78×	2.78×	57.19	64.08
CenterPoint-Ada3D ( <i>voxel=0.075</i> )	3.34×	3.96×	58.62	65.68

Outperform model-level compression methods

# Qualitative Results

- **Visualization of 2D & 3D heatmap**

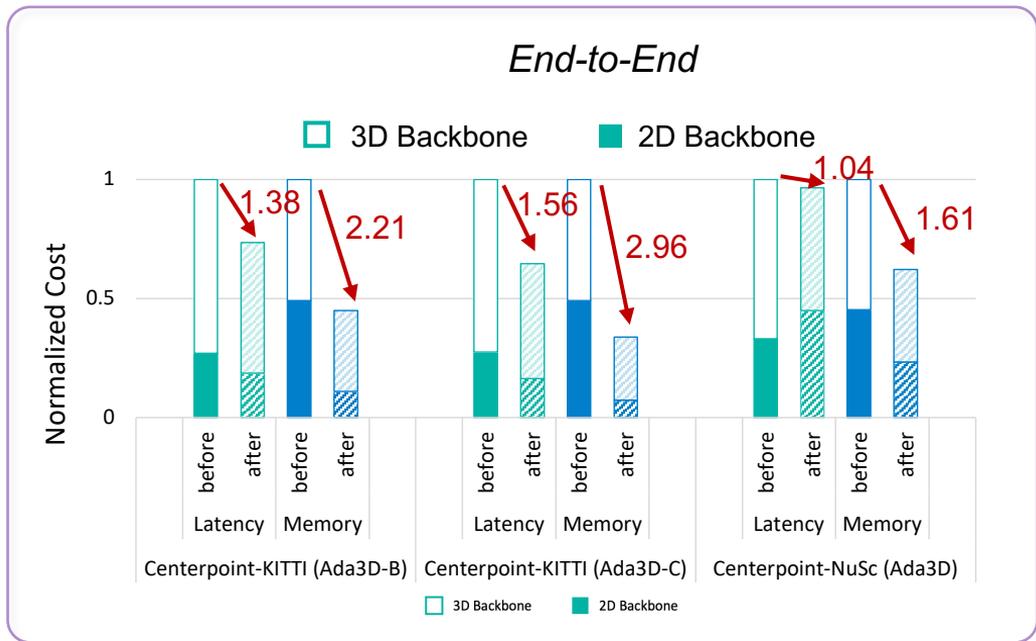
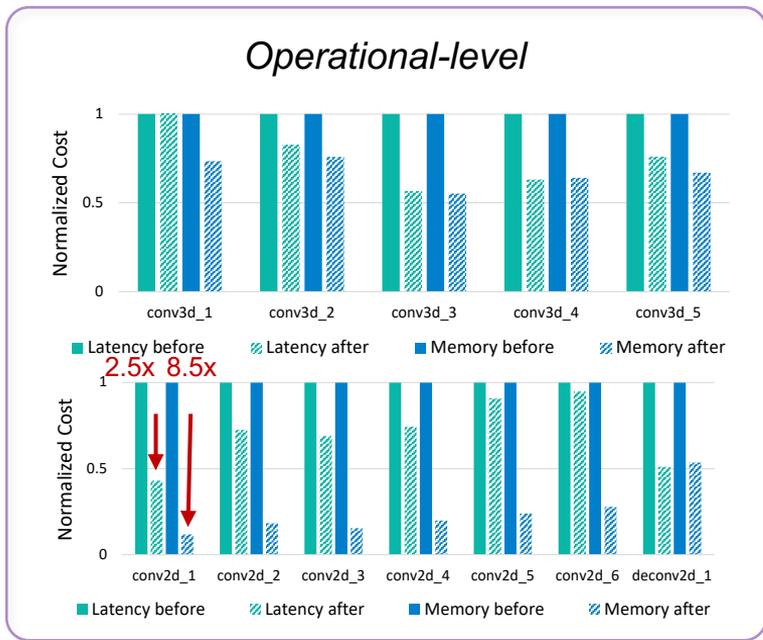
- The predictor identifies the voxels/pixels inside the bounding box
- The spatial filtering avoids dropping inbox data



# Hardware Experiments

## • Hardware Profiling on RTX3090 GPU

- Ada3D-B achieves **1.38x** latency and **2.21x** GPU (RTX 3090) peak memory optimization



# Summary

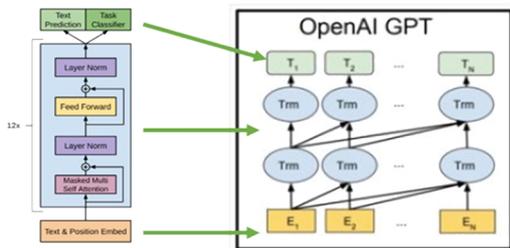
## Real-time Requirement



E.g., Video Conference **~30FPS**  
Mobile Phone / PC (1~10 TOPS)

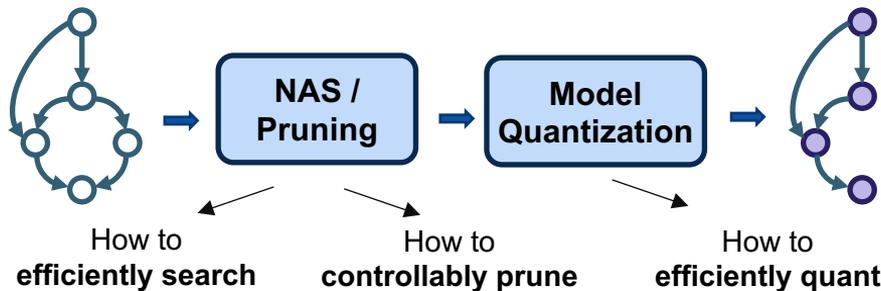


## High Latency

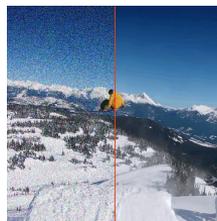


E.g., GPT-3: 175B, 664TOPs, **32s/seq**  
NVIDIA A100 (19.5 TFLOPS)

## Research Direction: Model Compression



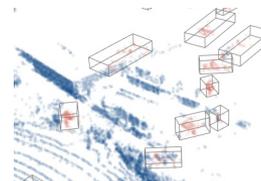
## Apply to practical scenarios



**Low-level Vision**  
**2x** peak mem. reduce



**Image Generation**  
**2~5x** speedup



**3D Perception**  
**1.38x** speedup & **2.21x**  
peak mem. reduce



清华大学电子工程系

Department of Electronic Engineering, Tsinghua University

# Thank You !

**Xuefei Ning** [foxdoraame@gmail.com](mailto:foxdoraame@gmail.com)

Department of Electronic Engineering, Tsinghua University

2023.03.16

Lab Leader: Prof. Yu Wang [yu-wang@tsinghua.edu.cn](mailto:yu-wang@tsinghua.edu.cn)

Thanks to our students: Zixuan Zhou, Tianchen Zhao, Shiyao Li, Xiangsheng Shi,

Lidong Guo, Junbo Zhao, Enshu Liu

