# Dynamic Networks

PD Dr. Haojin Yang

Multimedia and Machine Learning Group

Hasso Plattner Institute

**Design IT.**
**Create Knowledge.**
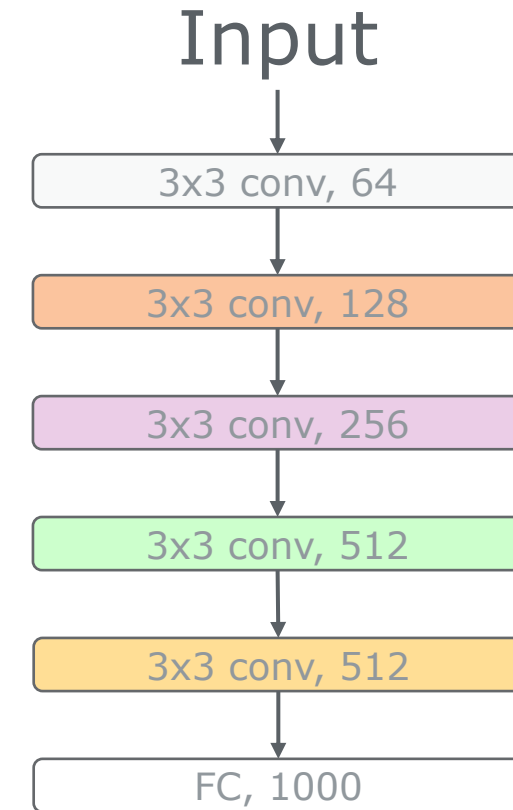
www.hpi.de

# Dynamic Networks

**Why Dynamic?**

Workflow of a static network

- Training

  - Pre-defined network architecture

  - Weight initialization

  - Optimization using training data

- Inference

  - **Fixed** network **architecture** and **weight**
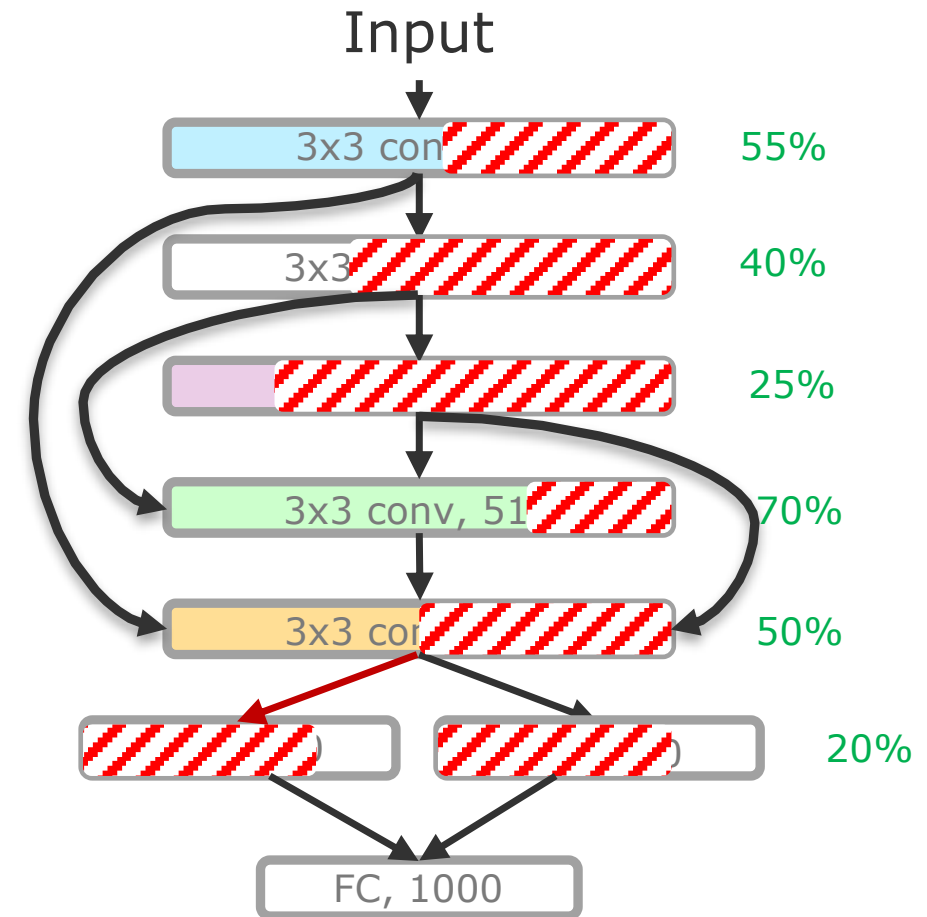
  - Forward pass for arbitrary input samples

Input

| 3x3 conv, 64 |
| 3x3 conv, 128 |
| 3x3 conv, 256 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| FC, 1000 |

# Dynamic Networks

**Why Dynamic?**

Workflow of a dynamic network

- Training
    - Learn a specific network structure for each sample
    - Joint training of neural network and decision-making mechanism
- Inference
    - Instance-wise dynamic network structure
    - The decision-making mechanism predicts the structure based on each input sample.
- Type:
    - Dynamic width
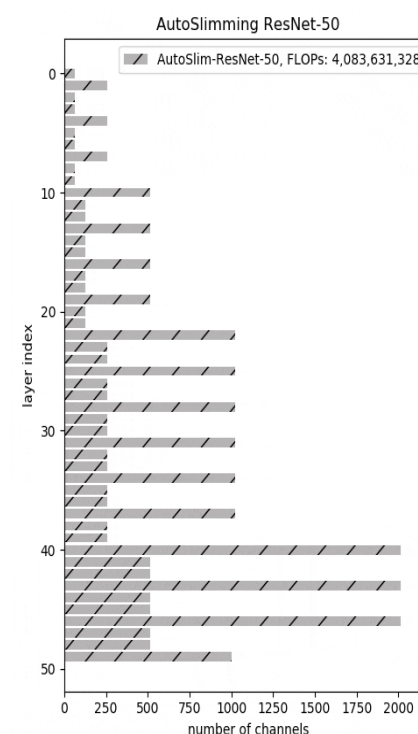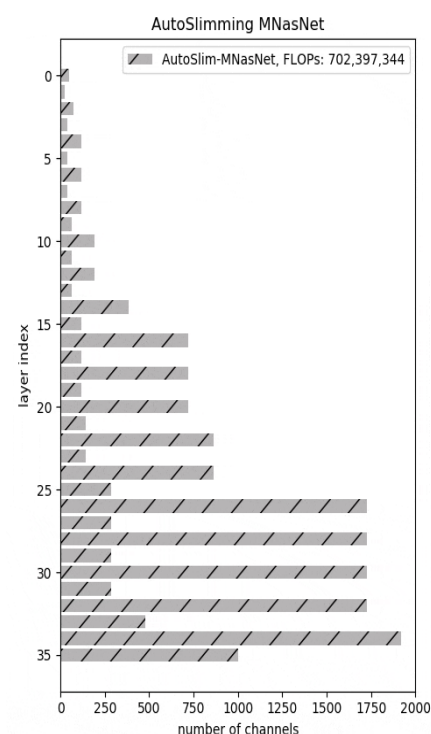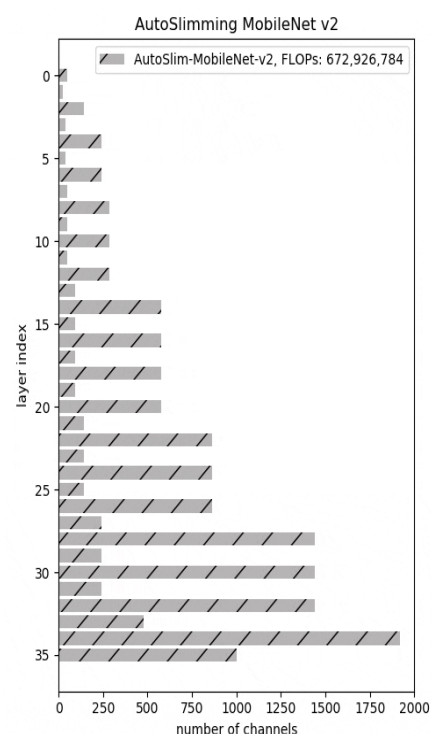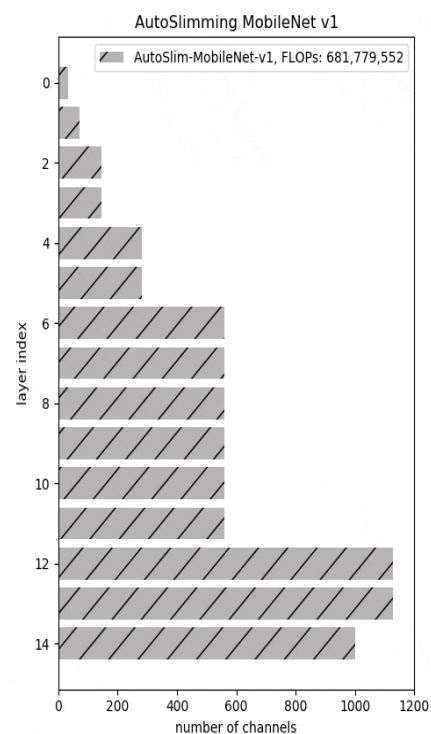    - Dynamic depth
    - Dynamic path

# Advantages

- Efficiency: Strategically allocate computations on demand at test time

- Representation power: Enlarged parameter space and improved representation power

- Adaptiveness: Desired trade-off between accuracy and efficiency on the fly

- Compatibility: Compatible with most advanced techniques in deep learning

- Generality: Can be applied seamlessly to a wide range of applications

- Interpretability: It is believed that the brains process information in a dynamic way

# Dynamic width

1.0×   0.75×   0.5×   0.25×

AutoSlimming MobileNet v1

AutoSlim-MobileNet-v1, FLOPs: 681,779,552

AutoSlimming MobileNet v2

AutoSlim-MobileNet-v2, FLOPs: 672,926,784

AutoSlimming MNasNet

AutoSlim-MNasNet, FLOPs: 702,397,344

AutoSlimming ResNet-50

AutoSlim-ResNet-50, FLOPs: 4,083,631,328

*Yu, Jiahui, and Thomas Huang. "Autoslim: Towards one-shot architecture search for channel numbers." arXiv preprint arXiv:1903.11728, 2019*

# Dynamic depth

- Early Exits



| 3x3 conv, 64 |
| 3x3 conv, 64,pool |
| 3x3 conv, 128 |
| 3x3 conv, 128,pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256,pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512,pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512,pool |
| FC, 4096 |
| FC, 4096 |
| FC, 1000 |

**Inference**

Confident? Yes: exit; No: go to the next exit

**Classifier 4**

Confident? Yes: exit; No: go to the next exit

**Classifier 3**

Confident? Yes: exit;
No: go to the final exit

**Classifier 2**

Joint
Training

**Final exit**

**Classifier 1**

$$Loss = \lambda_1 L_{cls1} + \lambda_2 L_{cls2} + \lambda_3 L_{cls3} + \lambda_4 L_{cls4}$$
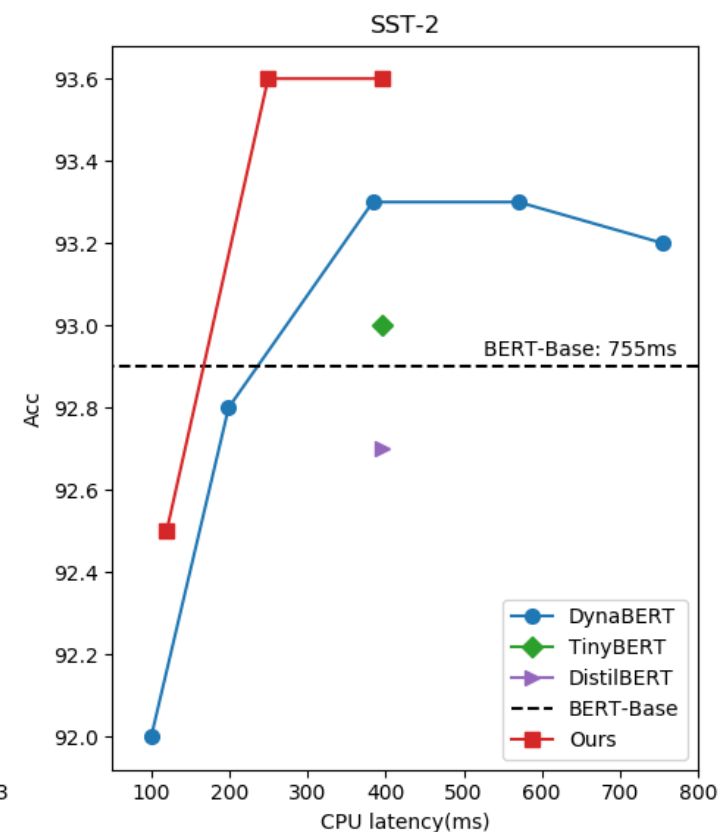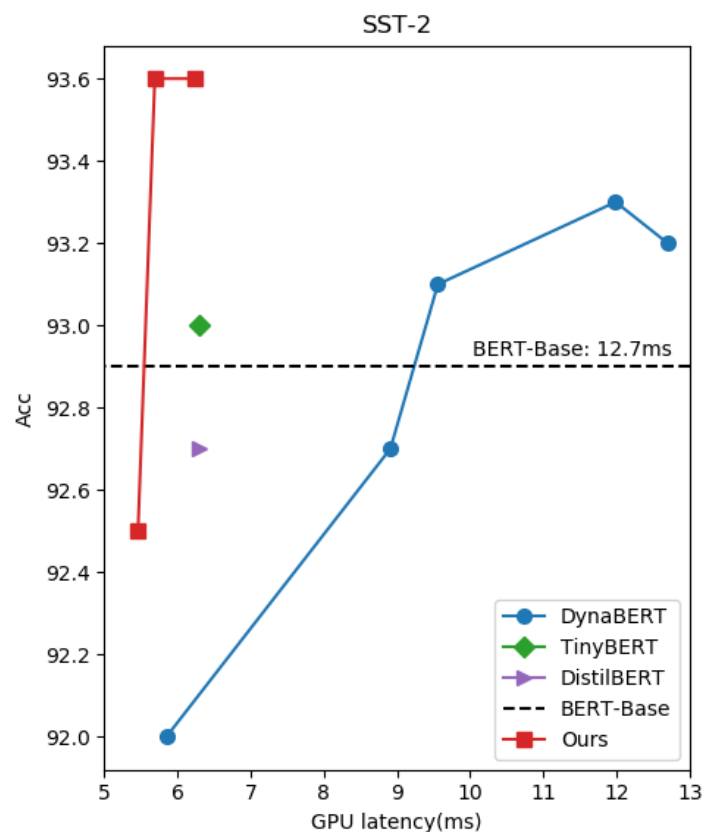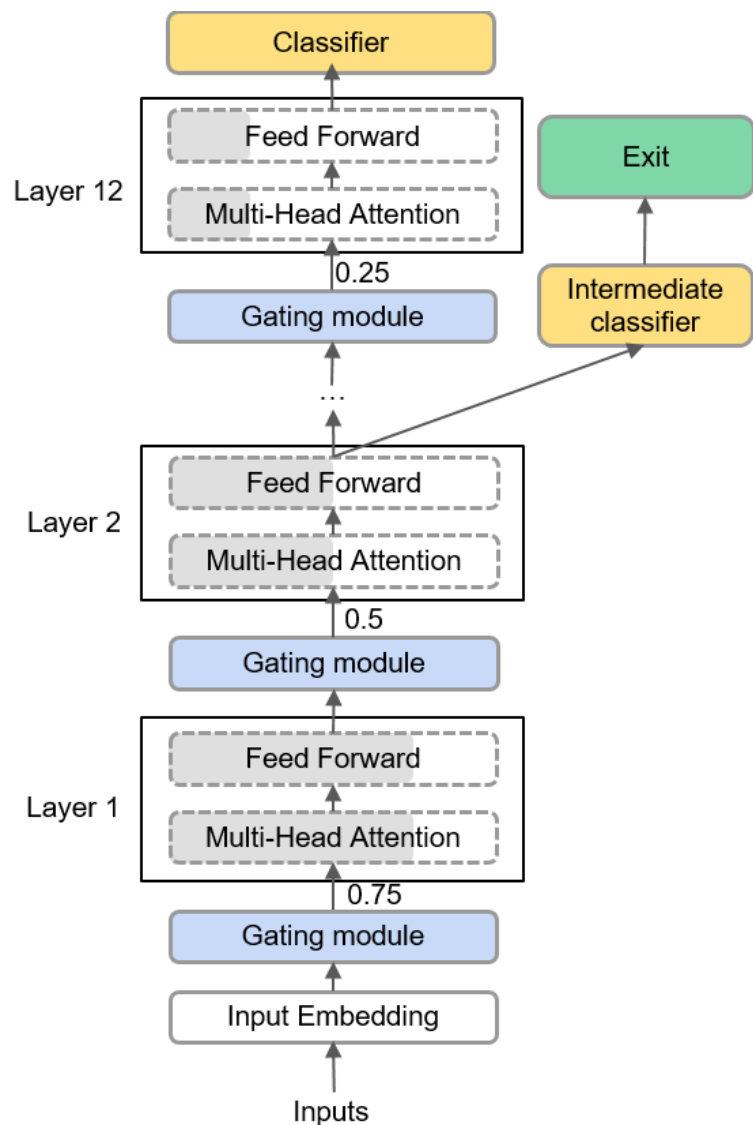
# Case study: Dynamic BERT

# Summary

- Differences between static and dynamic networks

- Advantages of dynamic networks

- Case study

  - Dynamic width

  - Dynamic depth

- Dynamic BERT example

- Upcoming: Challenges and solutions for efficient LLM deployment