# Efficient Solution to Large-scale Image Classification

Presenter: Chenhao Lin

Team: BigVideo

16/06/2019

# Team: BigVideo

Team Member:

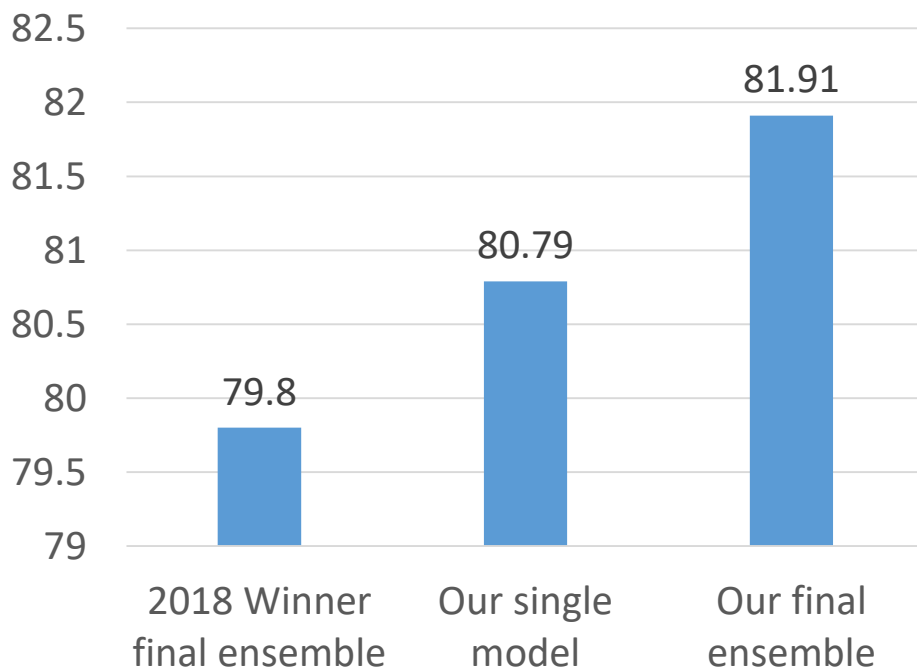Huabin Zheng     Litong Feng     Yuming Chen     Weirong Chen
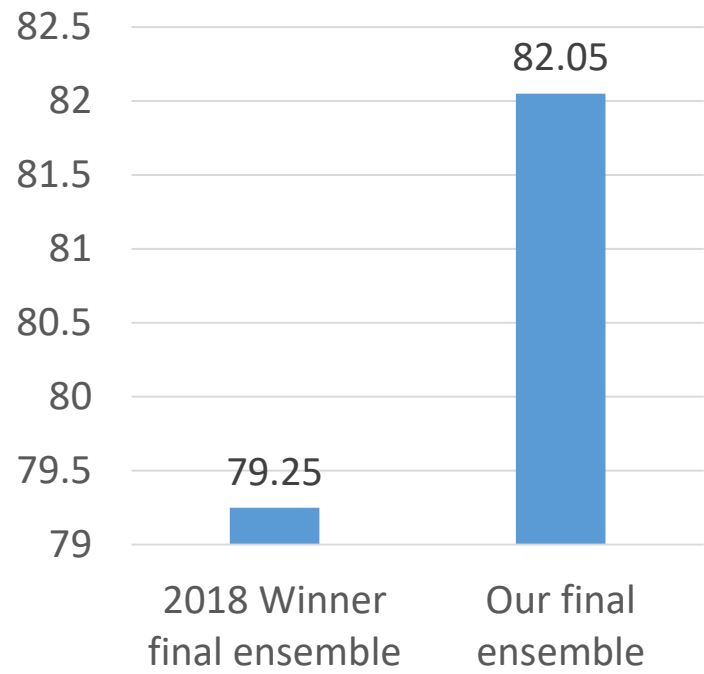
Zhe Huang          Zhanbo Sun          Wayne Zhang

# Results



Validation Top5

| | |
|---|---|
| 2018 Winner final ensemble | 79.8 |
| Our single model | 80.79 |
| Our final ensemble | 81.91 |

Test Top5

| | |
|---|---|
| 2018 Winner final ensemble | 79.25 |
| Our final ensemble | 82.05 |

# Overview

**Challenge**:

Limited GPU resources

VS

Large-scale data

Idea Validation

Many-model Ensemble

Pipeline:

**Model Selection**

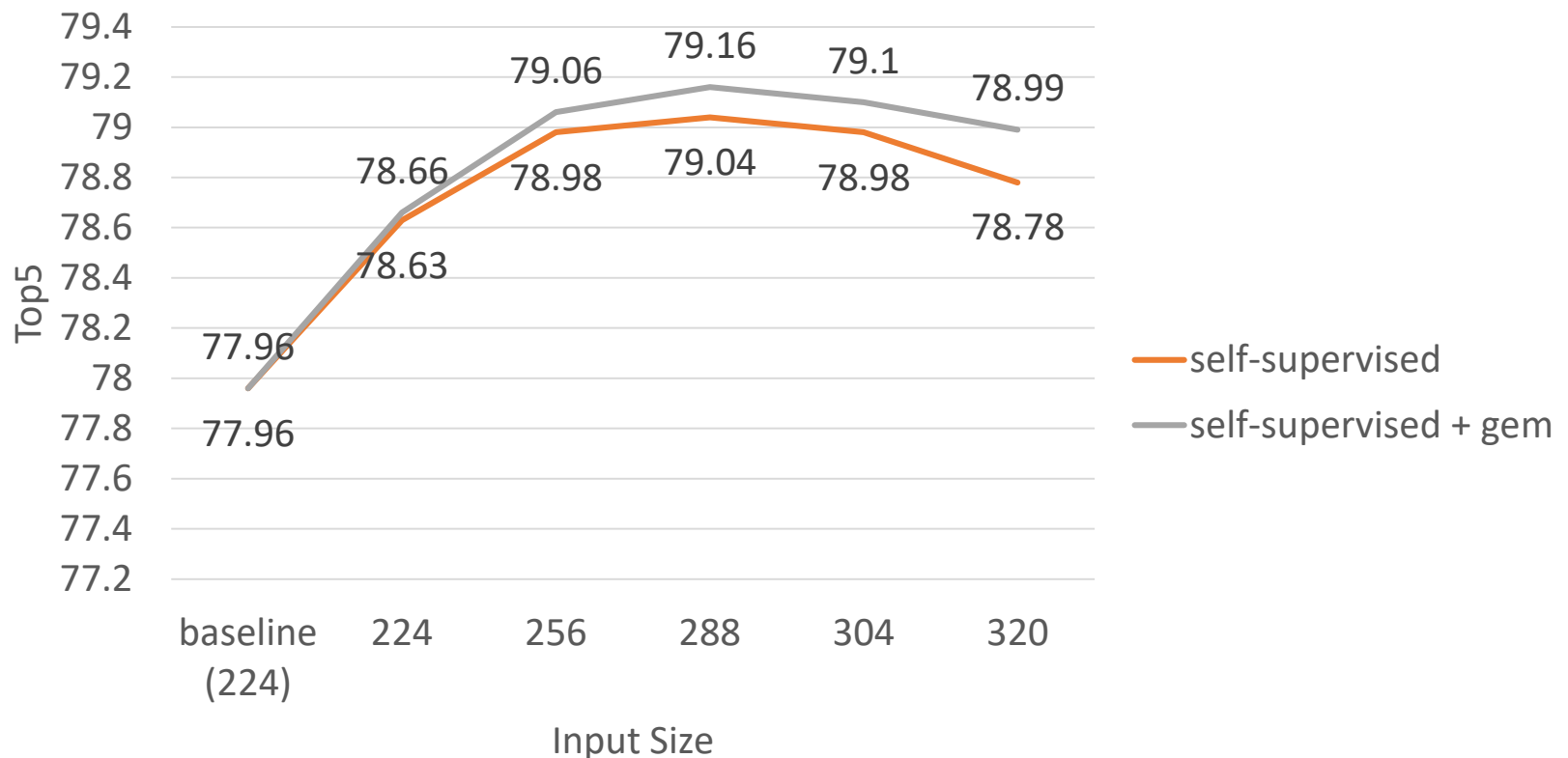Efficient & Powerful Network Architectures

↓

**Model Training (Starting from our in-house image classification tool)**

ImageNet-style Training

→

**Fine-tuning**

Large Input Size

Self-Supervised Loss

Using Description Text

→

**Final Ensemble of Three Models**

Multi-Crop Testing

4

# Efficient & Powerful Networks

| Network (Input Size) | ImageNet Top1 | Estimated Training Time on WebVision* |
|---|---|---|
| NASNet-A (331) | 82.70 | 64 GPUs 67 days |
| PNASNet-5 (331) | 82.90 | 64 GPUs 61 days |
| SENet154 (224) | 81.32 | 64 GPUs 18 days |
| **ResNeXt152 variant (224) (Our Primary Model)** | **81.53** | **64 GPUs 12 days** |
| Inception-ResNet-v2 (299) | 80.10 | 64 GPUs 12 days |
| DPN98(224) | 79.80 | 64 GPUs 11 days |
| SEResNet152(224) | 78.43 | 64 GPUs 9 days |

*Estimated training time for Webvision 150 epochs on TITANXp

# Fine-tuning with Expanded Input Size
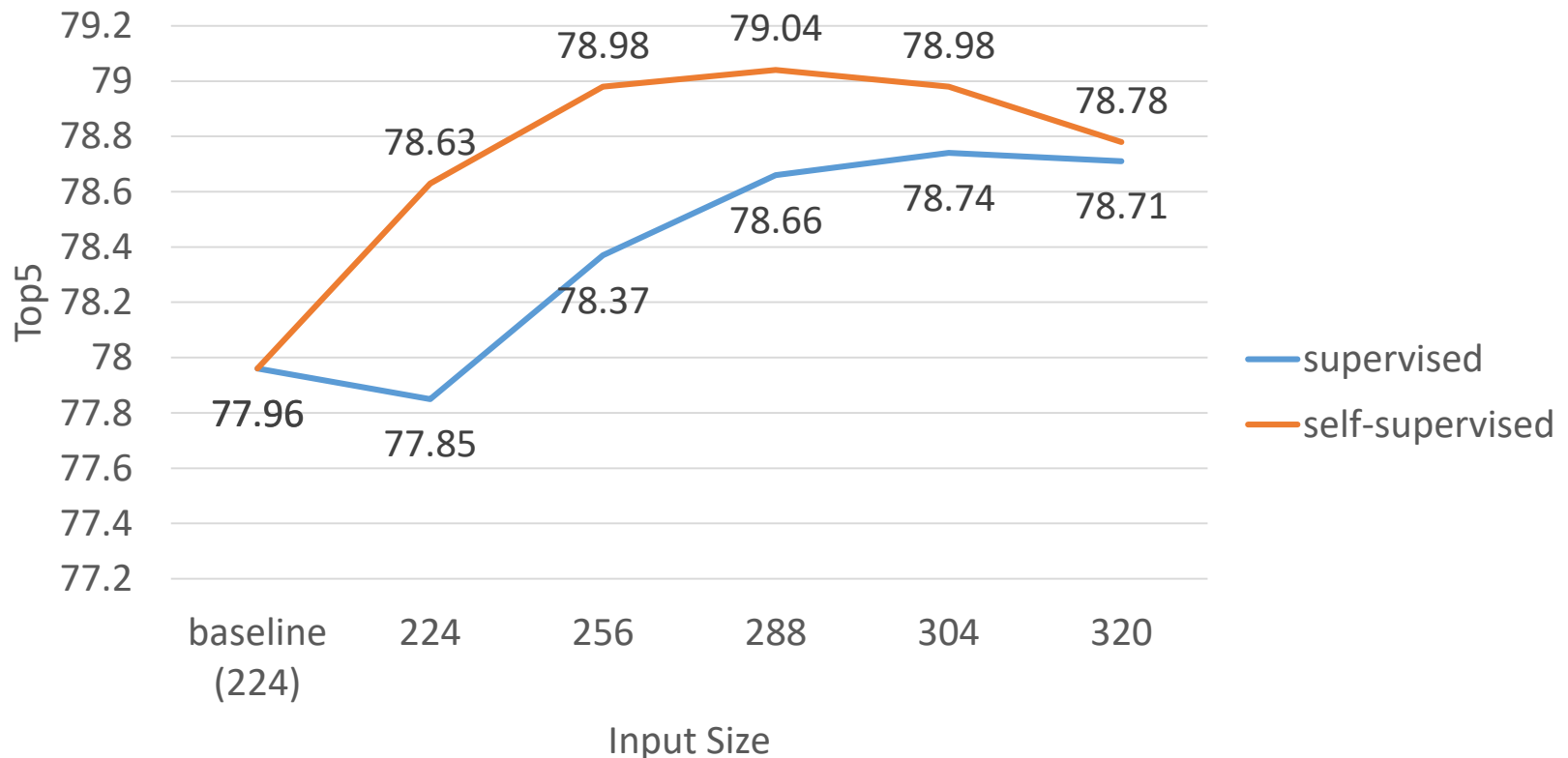
❑ Experience from ImageNet:
   ❑ Larger input size performs better.
   ❑ Due to limited resources, we fine-tune with large input sizes only.
❑ Generalized-Mean (GeM) pooling [1] adapts with large inputs better than global average pooling.

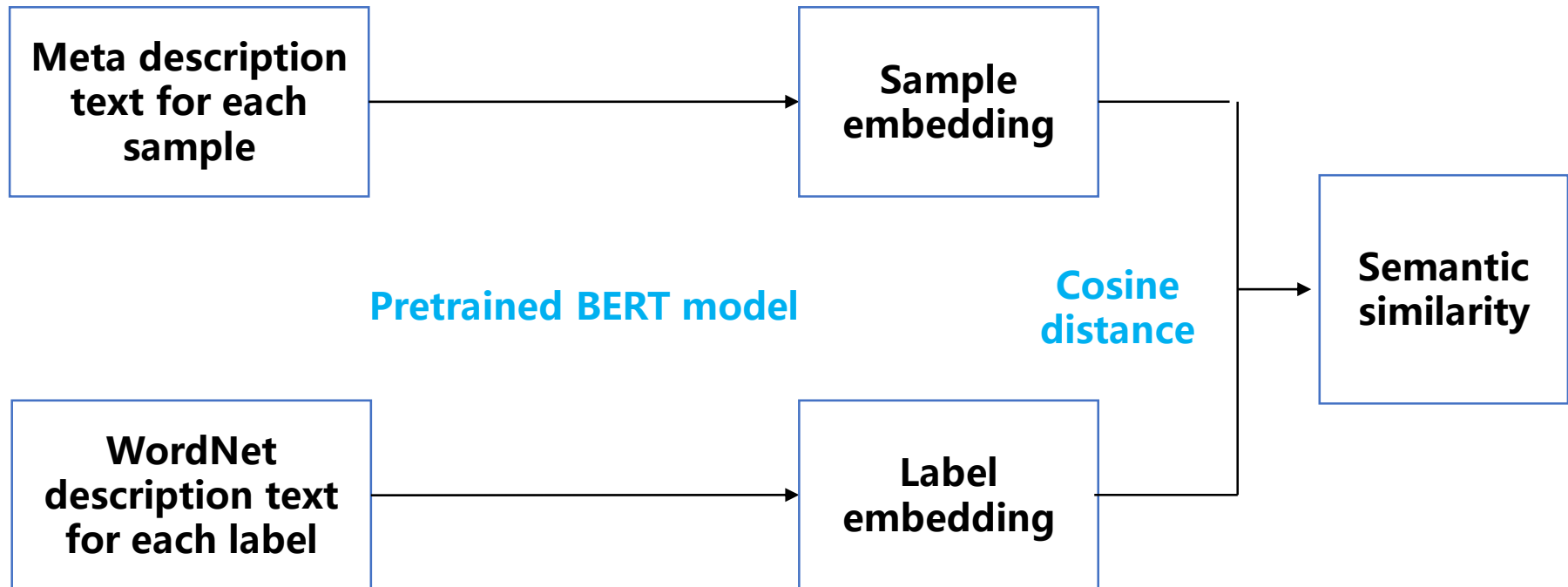[1] Berman, Maxim, et al. "MultiGrain: a unified image embedding for classes and instances."

# On-the-fly Self-supervised Loss

❑ After supervised training converges, pseudo labels from network itself are more reliable than noisy ground-truth labels.

# Using Description Text

❑ Select samples by semantic similarity between embeddings of sample description text and label description text.



**Meta description text for each sample** → **Sample embedding**

**Pretrained BERT model**

**WordNet description text for each label** → **Label embedding**

**Cosine distance**

**Semantic similarity**

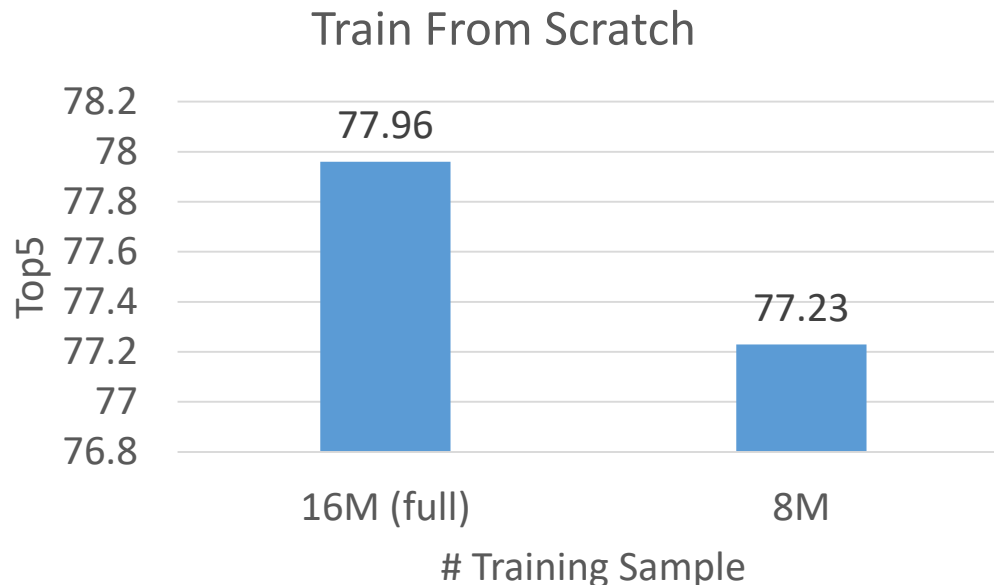# Using Description Text

Tag: Yardbird



random selection

top 10

middle 10
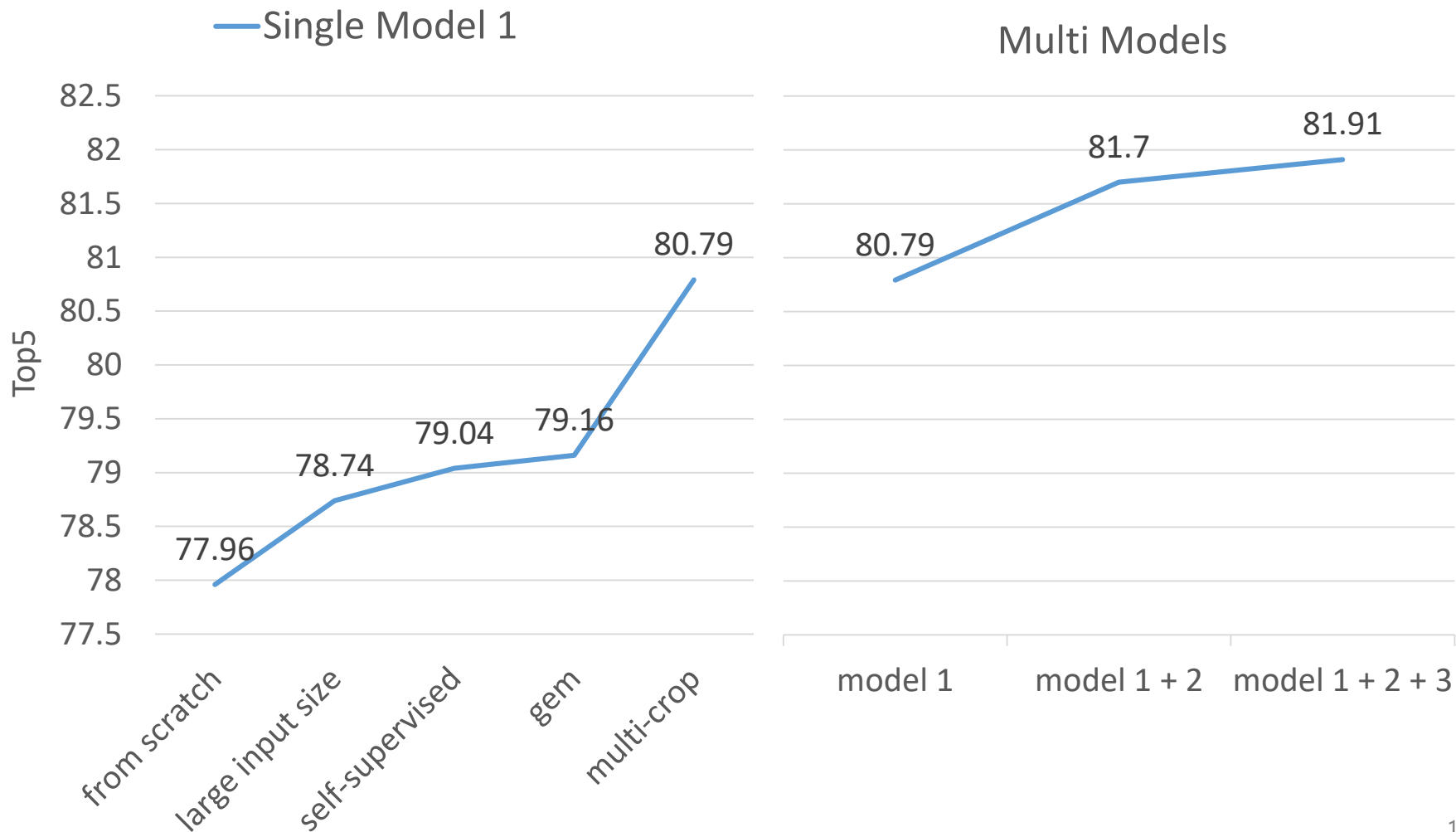
bottom 10

# Using Description Text

❑ Despite of visually appealing selection, we found training from scratch with the selected partial training set did not perform as well as with the full training set.
❑ Nevertheless, partial-set model contributes to the final ensemble's performance.

Train From Scratch

# Ensemble

# Take-home Message

❑ Fundamental improvements of image classification bring large gains.
- Efficient network with large capacity
- Expanded input size + GeM pooling
- On-the-fly self-supervised loss

❑ Side information may bring gains, however we did not have enough time and GPUs to explore them.
- Description text based sample selection using BERT

❑ De-noising tricks are hard to tune well.
- GHM
- Focal loss

# BigVideo Research Team of SenseTime

Dedicated to research on deep understanding of Internet photos & videos

- Holistic Semantic Understanding
  - People, Scene, Action, Event

- Big Data
  - 1 billion Images/Frames processed per day

- High Accuracy
  - 90% recall @ 1 / 1,000,000 FAR

- High Performance
  - 3000 QPS single GPU

50+ Researchers, 8 PhDs, 100+ Publications

# Thank You!