**Stable Diffusion/DALL-E/Imagen**

- High-quality realistic images



"cat"

Rombach et al., **"High-Resolution Image Synthesis with Latent Diffusion Models "**, CVPR, 2022
Ramesh et al. "**Zero-shot text-to-image generation**", ICML 2021
Saharia et al., "**Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding**", NeurIPS, 2022

## Stable Diffusion/DALL-E/Imagen

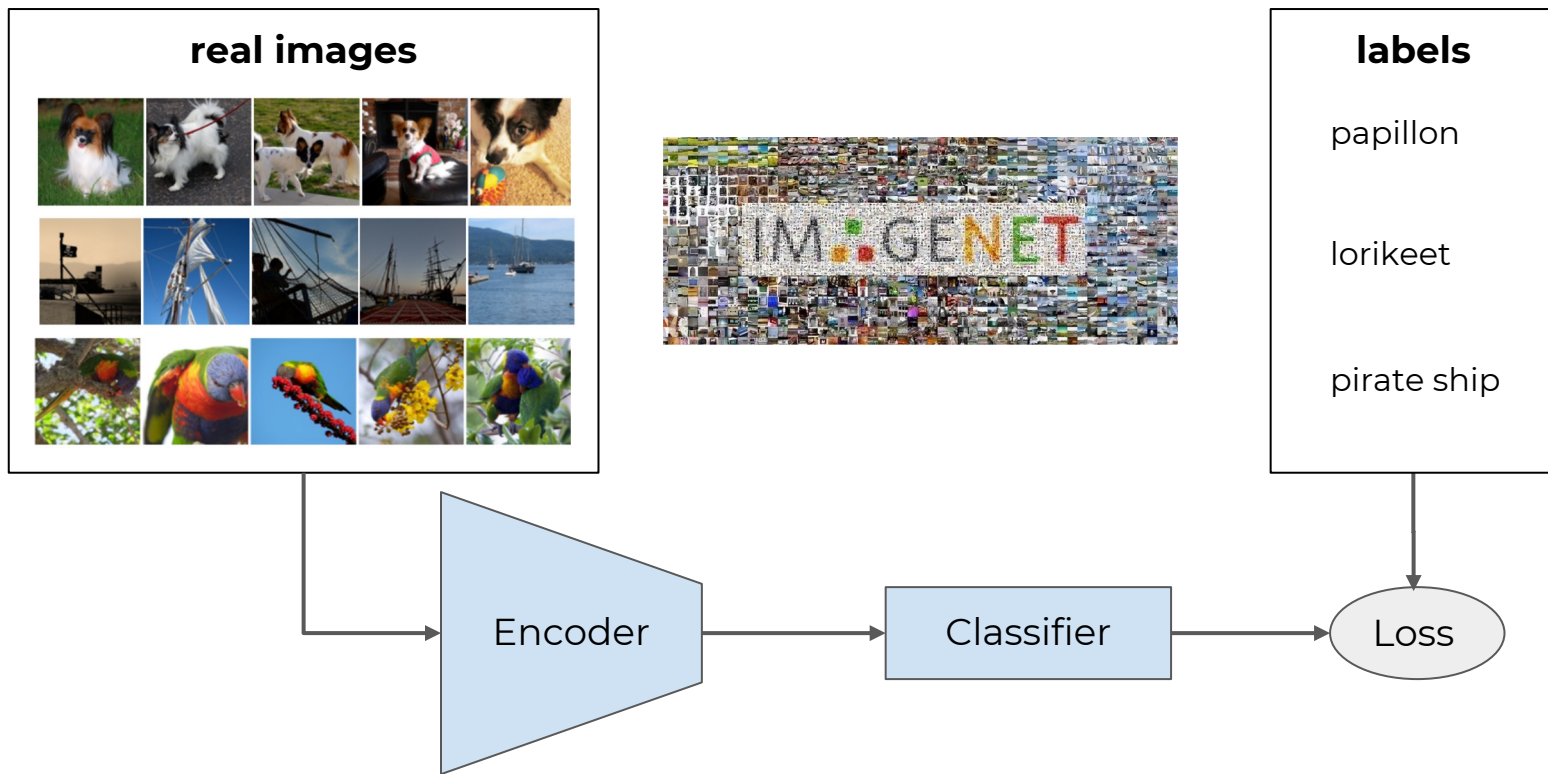- High-quality realistic images



"cat"

***Do we still need real images
for learning visual representations?***

Rombach et al., **"High-Resolution Image Synthesis with Latent Diffusion Models "**, CVPR, 2022
Ramesh et al. "**Zero-shot text-to-image generation**", ICML 2021
Saharia et al., "**Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding**", NeurIPS, 2022

**real images**

**labels**

papillon

lorikeet

pirate ship

Encoder

Classifier

Loss

Deng et al. "**Imagenet: A large-scale hierarchical image database**", CVPR 2009
Russakovsky et al., **"Imagenet large scale visual recognition challenge**", IJCV, 2015

# Can **ImageNet-1K** be replaced by synthetic images?



synthetic images

labels

papillon

lorikeet

pirate ship

Encoder → Classifier → Loss

Deng et al. "**Imagenet: A large-scale hierarchical image database**", CVPR 2009
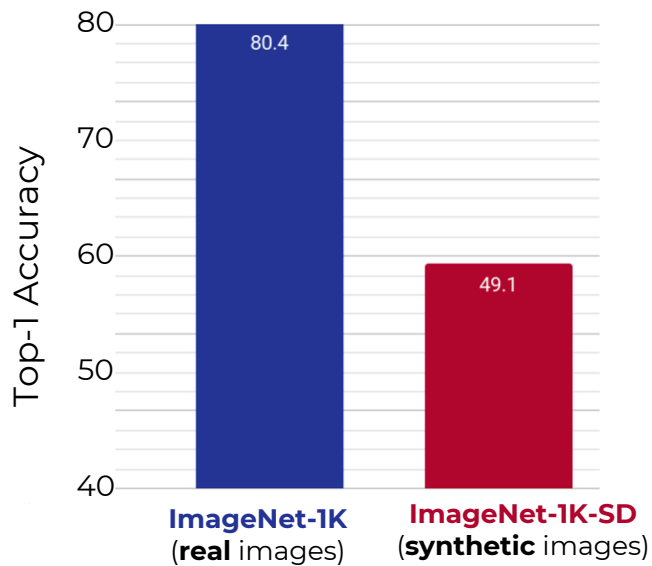Russakovsky et al., "**Imagenet large scale visual recognition challenge**", IJCV, 2015

# Training image classifiers on **ImageNet-SD**
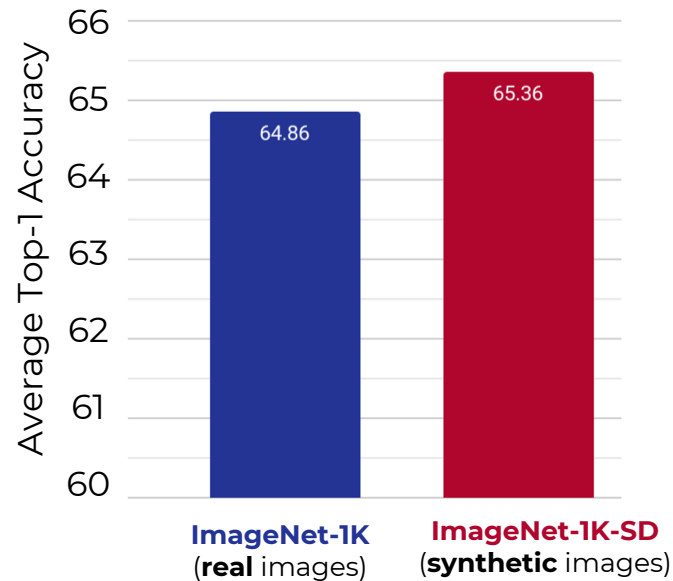
# Overview of the results

Performance on **ImageNet-1K** val. set
(*real* images)

Performance on **15 transfer datasets**
(*real* images)

Textual
Prompt

Stable
Diffusion

Synthetic Image

Rombach et al., **"High-Resolution Image Synthesis with Latent Diffusion Models** ", CVPR, 2022
Pretrained model available at **https://huggingface.co/CompVis/stable-diffusion-v1-4**

Textual
Prompt

Stable
Diffusion

Synthetic Image

prompt = class name

"papillon"



"lorikeet"



"pirate, pirate ship"

# Prompts for synthesizing ImageNet clones

Textual Prompt → Stable Diffusion → Synthetic Image

prompt = class name

"papillon"

"lorikeet"

"pirate, pirate ship"

Semantic errors

Lack of diversity

Domain issues

"papillon" class in ImageNet

"pirate, pirate ship" class in ImageNet

Textual
Prompt



Stable
Diffusion

Synthetic Image

prompt = class name, hypernym*

"papillon, <hypernym$^{papillon}$>"



"lorikeet, <hypernym$^{lorikeet}$>"



"pirate ship, <hypernym$^{pirate-ship}$>"



\* from **Wordnet** lexical database

Miller. et al. "**WordNet: A Lexical Database for English**". ACM-Comm 1995

# Tackling semantic & domain issues

Textual Prompt → Stable Diffusion → Synthetic Image

prompt = class name, hypernym*

"papillon, <hypernym$^{papillon}$>"



"lorikeet, <hypernym$^{lorikeet}$>"



"pirate ship, <hypernym$^{pirate-ship}$>"



prompt = class name, description*

"papillon, <description$^{papillon}$>"



"lorikeet, <description$^{lorikeet}$>"



"pirate ship, <description$^{pirate-ship}$>"



* from **Wordnet** lexical database

Miller. et al. "**WordNet: A Lexical Database for English**". ACM-Comm 1995

Textual
Prompt

Stable
Diffusion

Synthetic Image

prompt = class name, hypernym inside background**

"papillon, <hypernym$^{papillon}$>
inside <background>"



"lorikeet, <hypernym$^{lorikeet}$>
inside <background>"



"pirate ship, <hypernym$^{pirate-ship}$>
inside <background>"



**  from **Places 365** dataset

Zhou et al. "**Places: A 10 million image database for scene recognition**." PAMI, 2017

Textual
Prompt

Stable
Diffusion

Synthetic Image

prompt = class name, hypernym inside background**

"papillon, <hypernym$^{papillon}$>
inside <background>"

"lorikeet, <hypernym$^{lorikeet}$>
inside <background>"

"pirate ship, <hypernym$^{pirate\text{-}ship}$>
inside <background>"



prompt = class name, description (+ reduce guidance scale)

"papillon, <description$^{papillon}$>"

"lorikeet, <description$^{lorikeet}$>"

"pirate ship, <description$^{pirate\text{-}ship}$>"



**  from **Places 365** dataset

Zhou et al. "**Places: A 10 million image database for scene recognition**." PAMI, 2017
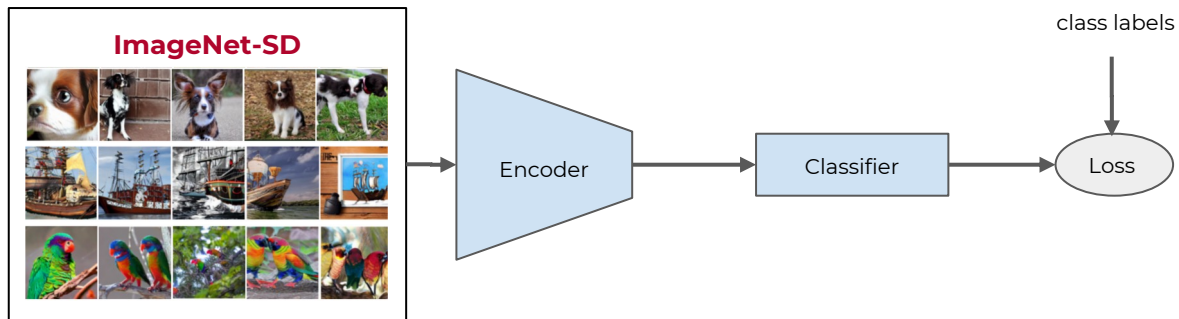
ImageNet-SD

**ImageNet-SD** datasets:

Synthetic clones of different ImageNet subsets

- **ImageNet-100-SD**: 100 classes, 130k images

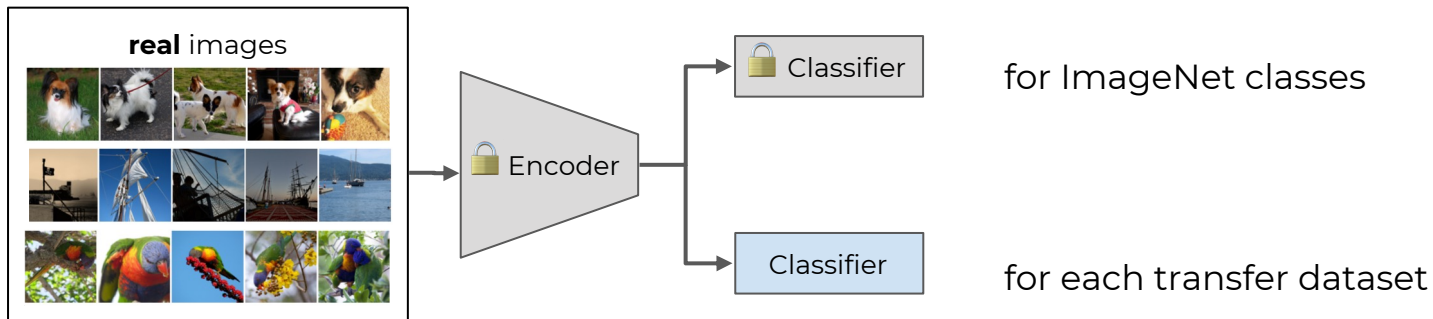- **ImageNet-1K-SD**: 1000 classes, 1.2M images

## Training with synthetic data



## Evaluation protocol

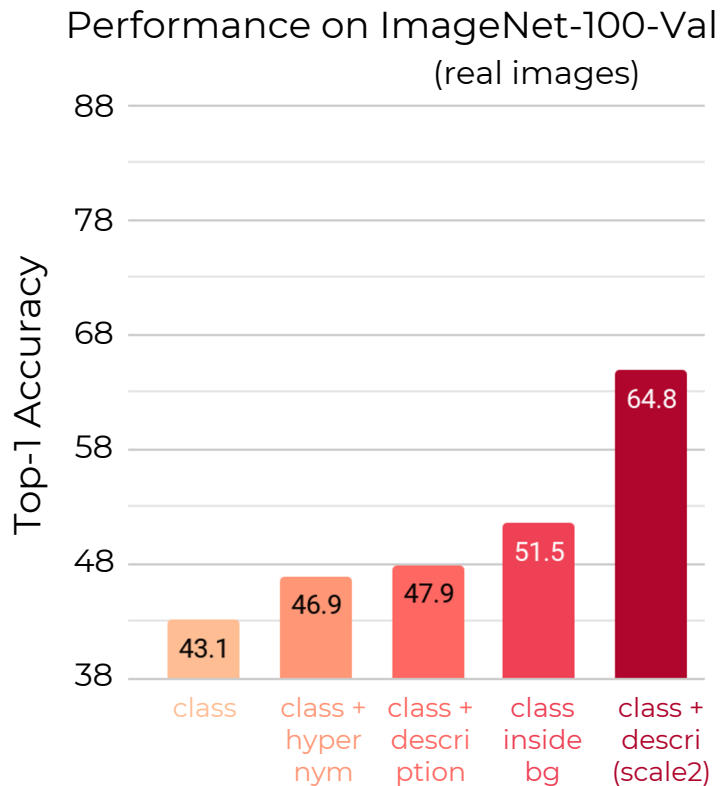Performance on ImageNet-100-Val
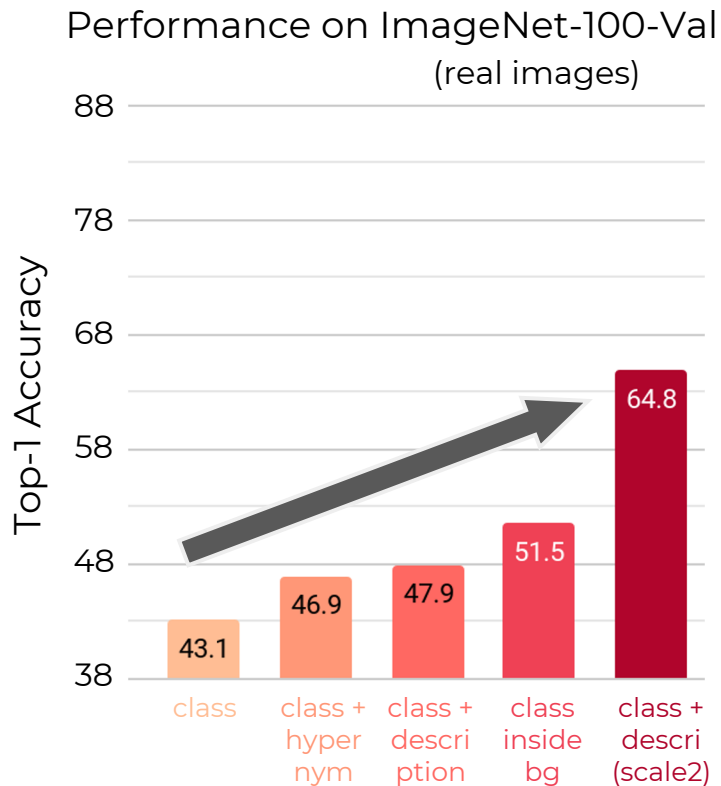(real images)

# ImageNet-100: Results for different prompts



Performance on ImageNet-100-Val
(real images)

**Observations:**

- Addressing semantic, domain and diversity issues leads to better performance

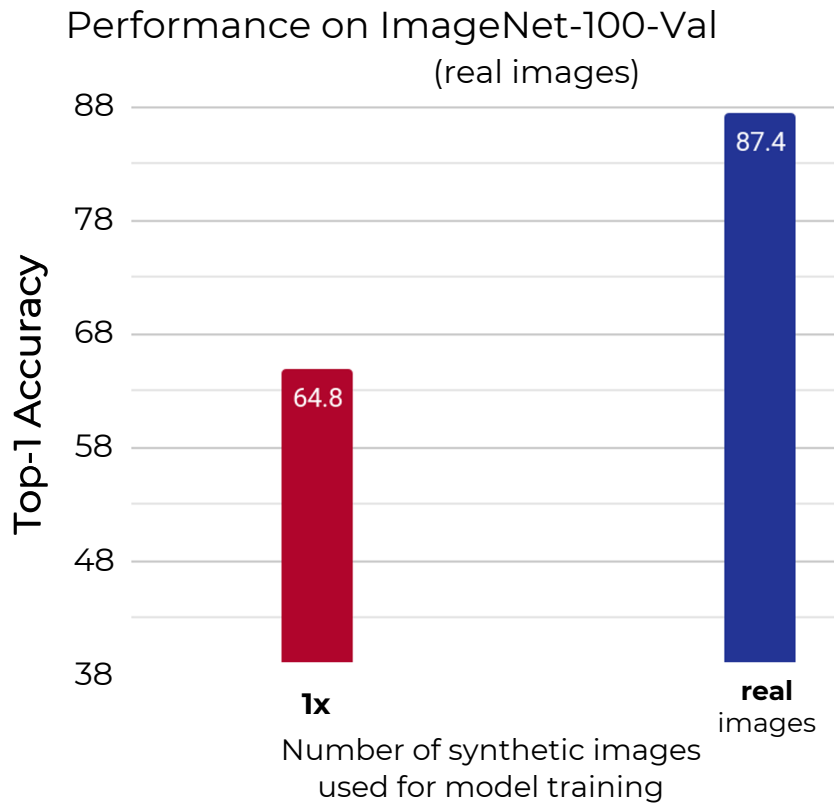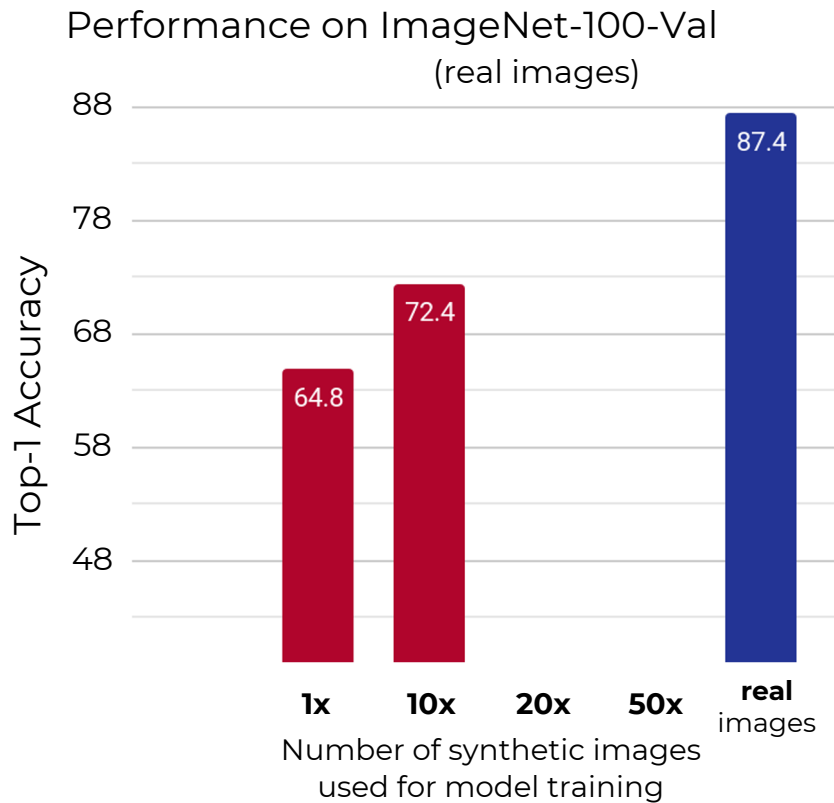Performance on ImageNet-100-Val
(real images)

**Observations:**

- Addressing semantic, domain and diversity issues leads to better performance

- Significant gap between the models trained on **real** vs. **synthetic** images for the training classes
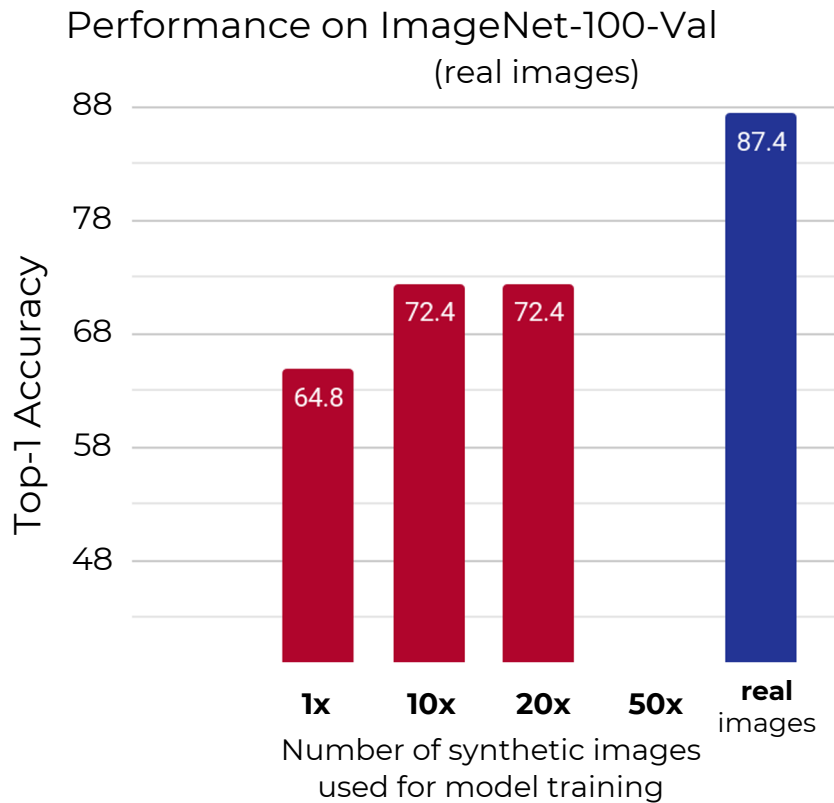
# ImageNet-100: Scaling the number of synthetic images

Performance on ImageNet-100-Val
(real images)

ImageNet-100: Scaling the number of synthetic images

# ImageNet-100: Results for transfer learning

## Performance on 10 transfer datasets
### (real images)



Average Top-1 Accuracy

| | 1x | real images |
|---|---|---|
| | 56.2 | 58.4 |

Number of synthetic images used for model training

Kornblith et al., "**Do better ImageNet models transfer better?**", CVPR, 2019
[Long-tail] Horn et al., "**The iNaturalist species classification and detection dataset**", CVPR, 2018

# ImageNet-100: Results for transfer learning



Performance on 10 transfer datasets

(real images)

**Observations:**

- Increasing the number of synthetic images leads to higher transfer learning performance

- Representations from the model trained **synthetic** images *outperform* the ones from **real** for transfer learning

Performance on 10 transfer datasets

(real images)

**Observations:**

- Increasing the number of synthetic images leads to higher transfer learning performance

- Representations from the model trained **synthetic** images *outperform* the ones from **real** for transfer learning



Performance on 10 transfer datasets
(real images)

# Results for transfer learning

**Observations:**

- Increasing the number of synthetic images leads to higher transfer learning performance

- Representations from the model trained **synthetic** images *outperform* the ones from **real** for transfer learning



Performance on 10 transfer datasets (real images)

# Results for transfer learning
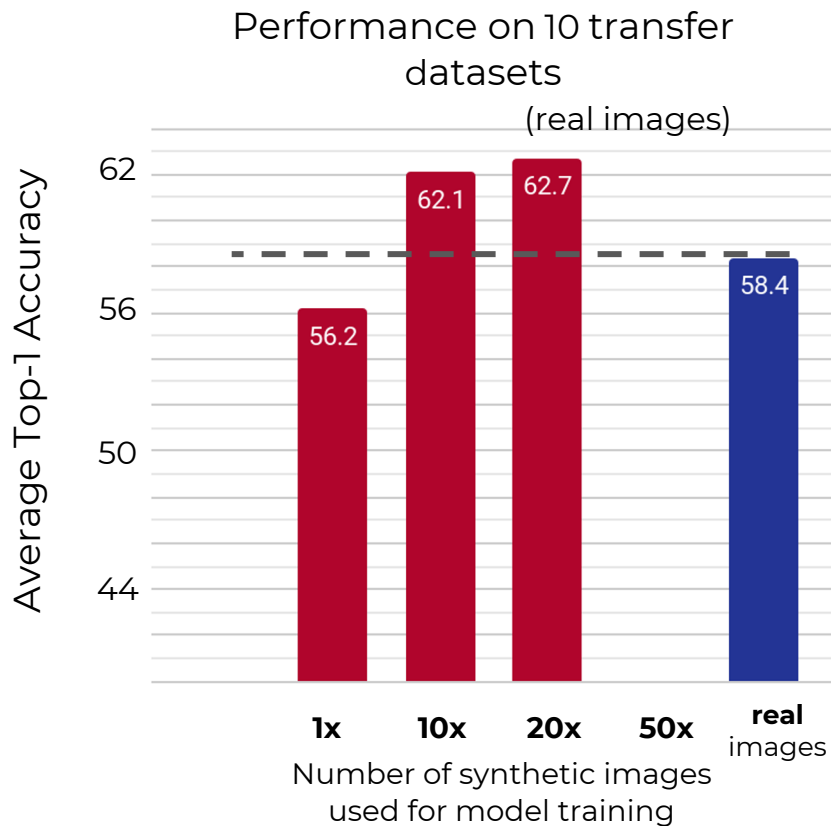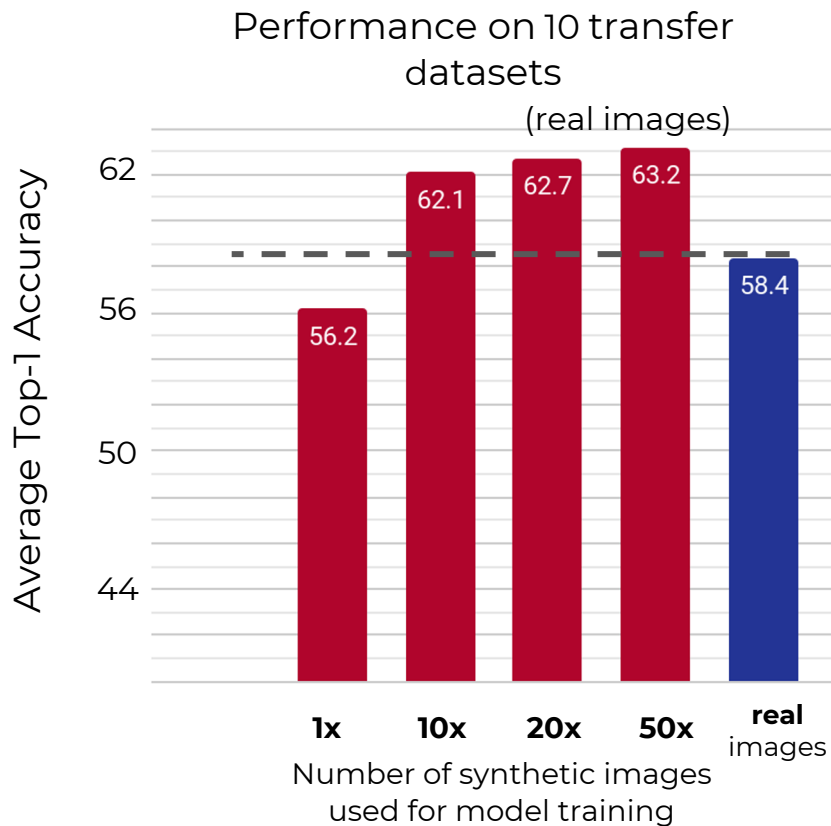
**Observations:**

- Increasing the number of synthetic images leads to higher transfer learning performance

- Representations from the model trained **synthetic** images *outperform* the ones from **real** for transfer learning



Performance on 10 transfer datasets
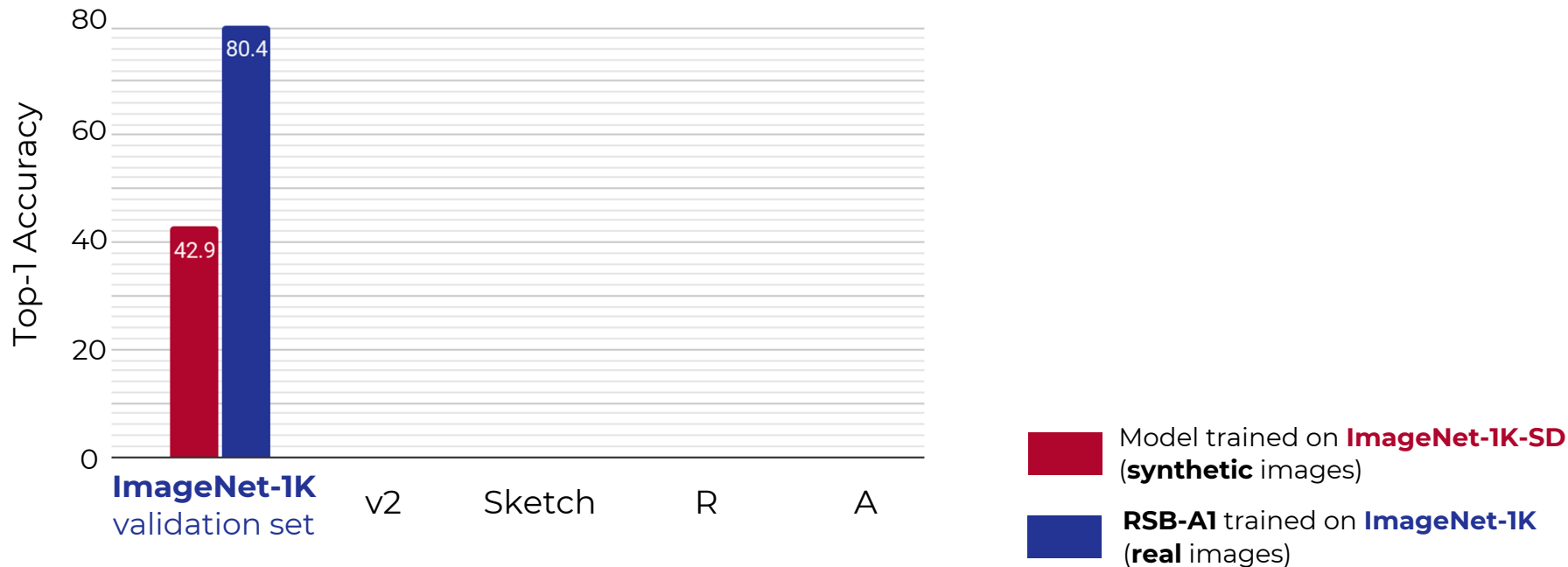(real images)

Average Top-1 Accuracy

| Label | Value |
|-------|-------|
| 1x | 56.2 |
| 10x | 62.1 |
| 20x | 62.7 |
| 50x | 63.2 |
| real images | 58.4 |

Number of synthetic images used for model training

# ImageNet-1K: Comparison to the state-of-the-art

Training with *the exact same number* of **real** and **synthetic** images per class

Top-1 Accuracy

80.4
42.9

ImageNet-1K validation set | v2 | Sketch | R | A

Model trained on **ImageNet-1K-SD** (**synthetic** images)

**RSB-A1** trained on **ImageNet-1K** (**real** images)

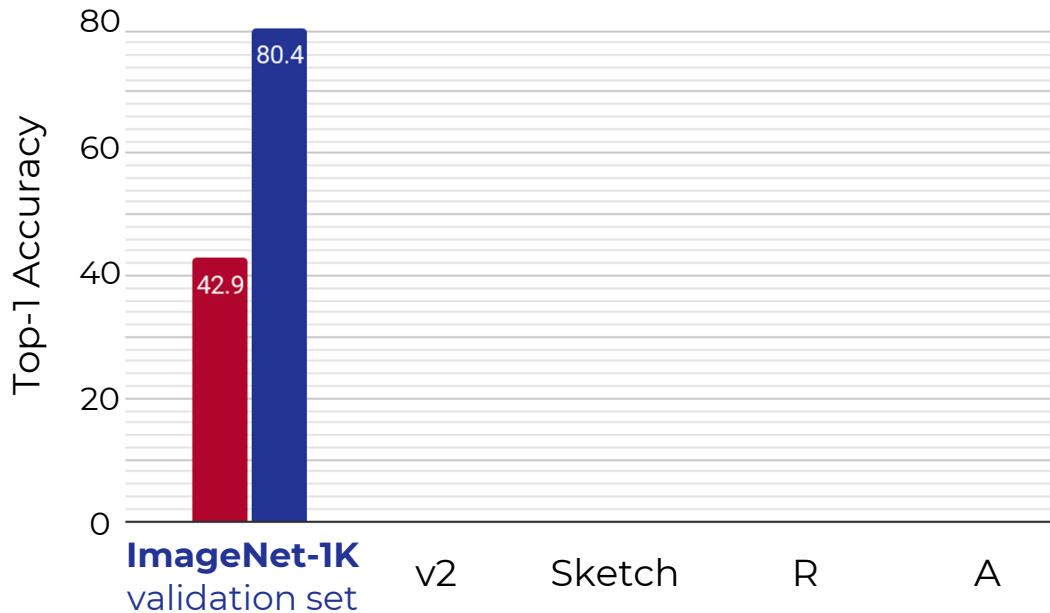[RSB-A1] Wightman et al., "**ResNet strikes back: An improved training procedure in timm**.", NeurIPSW, 2021

# ImageNet-1K: Comparison to the state-of-the-art

Training with *the exact same number* of **real** and **synthetic** images per class
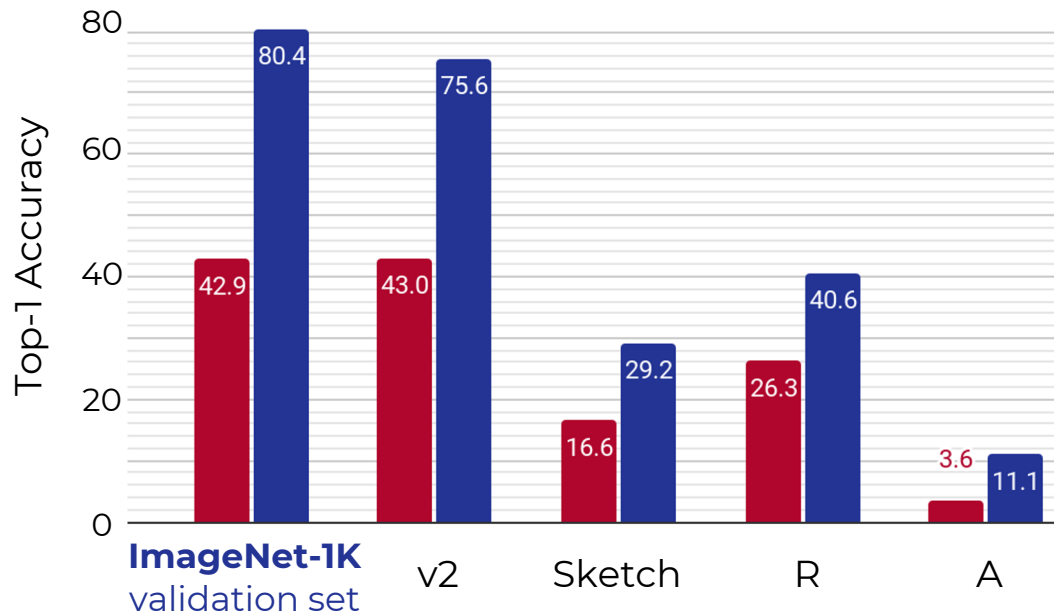


**Observations:**

- Significant gap between the models trained on **real** vs. **synthetic** images for the training classes

Model trained on **ImageNet-1K-SD** (**synthetic** images)

**RSB-A1** trained on **ImageNet-1K** (**real** images)

[RSB-A1] Wightman et al., "**ResNet strikes back: An improved training procedure in timm.**", NeurIPSW, 2021

# **ImageNet-1K:** Comparison to the state-of-the-art



**Observations:**

- Significant gap between the models trained on **real** vs. **synthetic** images for the training classes

- Relative gap is smaller for other variants especially ones with domain shifts

Model trained on **ImageNet-1K-SD** (**synthetic** images)

**RSB-A1** trained on **ImageNet-1K** (**real** images)

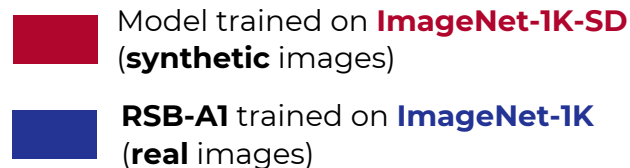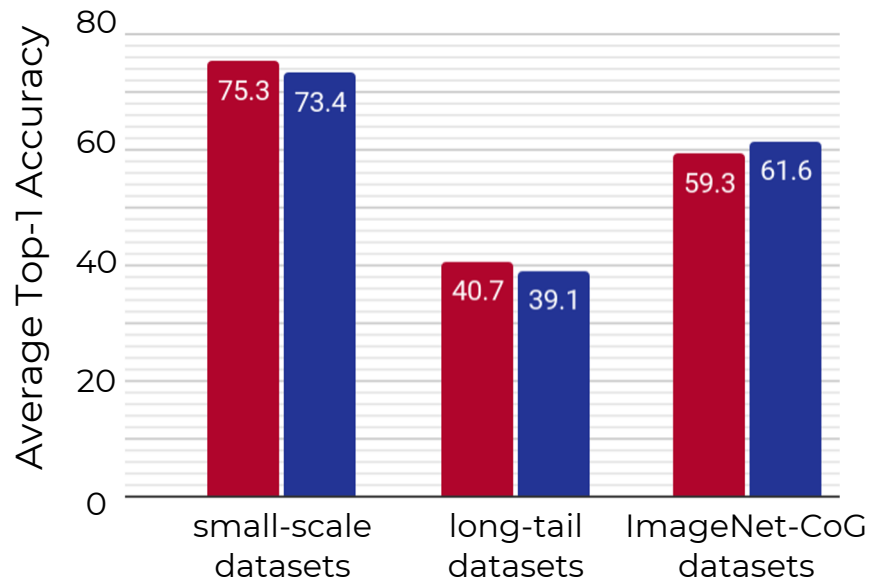[RSB-A1] Wightman et al., "**ResNet strikes back: An improved training procedure in timm**.", NeurIPSW, 2021

# **ImageNet-1K:** Comparison to the state-of-the-art



Performance on 15 transfer datasets
(real images)

[ImageNet-CoG] Sariyildiz et al., "**Concept Generalization in Visual Representation Learning**", ICCV, 2021
[Long-tail] Horn et al., "**The iNaturalist species classification and detection dataset**", CVPR, 2018
[Small-scale] Kornblith et al., "**Do better ImageNet models transfer better?**", CVPR, 2019

Model trained on **ImageNet-1K-SD**
(**synthetic** images)

**RSB-A1** trained on **ImageNet-1K**
(**real** images)

**Observations:**

- The model trained on **synthetic** images is *on-par or better* than the publicly available, state-of-the-art, **RSB-A1** model

- Synthesizing more images could lead to further gains

Performance on 15 transfer datasets
(real images)



[ImageNet-CoG] Sariyildiz et al., "**Concept Generalization in Visual Representation Learning**", ICCV, 2021
[Long-tail] Horn et al., "**The iNaturalist species classification and detection dataset**", CVPR, 2018
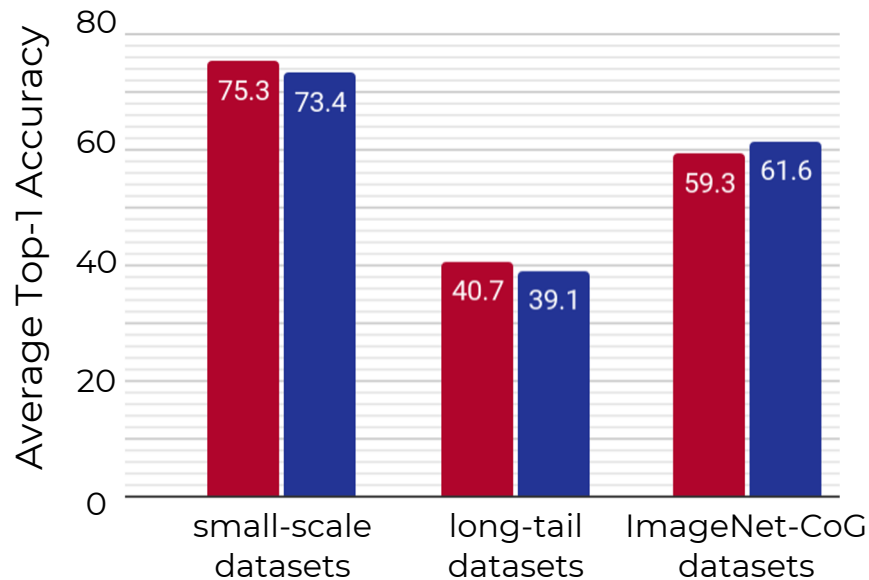[Small-scale] Kornblith et al., "**Do better ImageNet models transfer better?**", CVPR, 2019

Model trained on **ImageNet-1K-SD** (**synthetic** images)

**RSB-A1** trained on **ImageNet-1K** (**real** images)

What if we replace the **ImageNet** dataset
with **synthetic data** from **Stable Diffusion**?

**ImageNet-SD**:
Synthetic ImageNet clones
with Stable Diffusion images

**Result summary:**
- Decent but inferior performance on the ImageNet classes
- On-par or better performance than the state-of-the-art for transfer learning

**Bigger picture:**
- Image-free distillation of a generic text-to-image generation model
  into a visual encoder of arbitrary architecture, for solving a specific task

What if we replace the **ImageNet** dataset with **synthetic data** from **Stable Diffusion**?

**ImageNet-SD**:
Synthetic ImageNet clones with Stable Diffusion images

Come to our poster!
**TUE-PM-372**

Project page:
https://europe.naverlabs.com/imagenet-sd

p = name

p = name, hypernym

p = name, description

p = name inside background

p = name, description (scale=2)

real images