

Compatibility analysis of fine-tuned models for fine-grained text-based zero- shot image retrieval

András Schmelczer | MIR - 2022



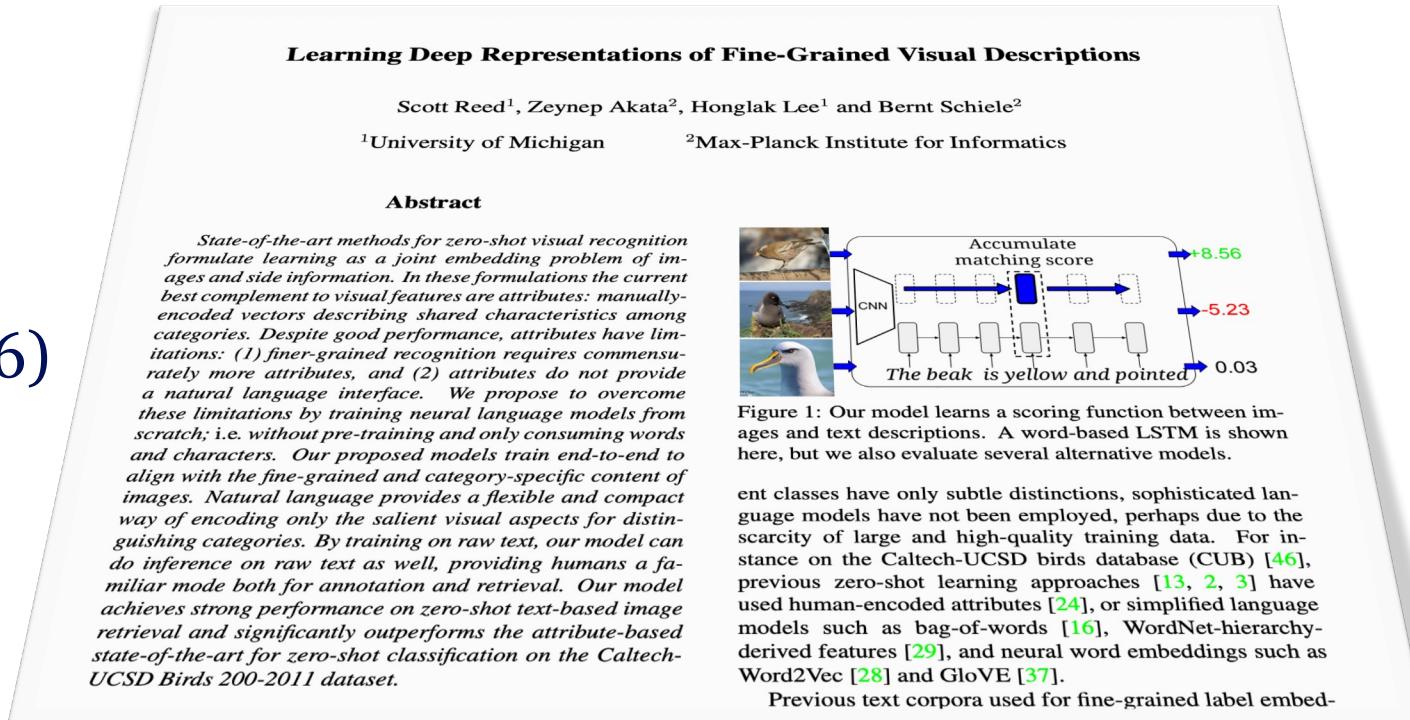
Universiteit
Leiden
The Netherlands

Discover the world at Leiden University

Motivation

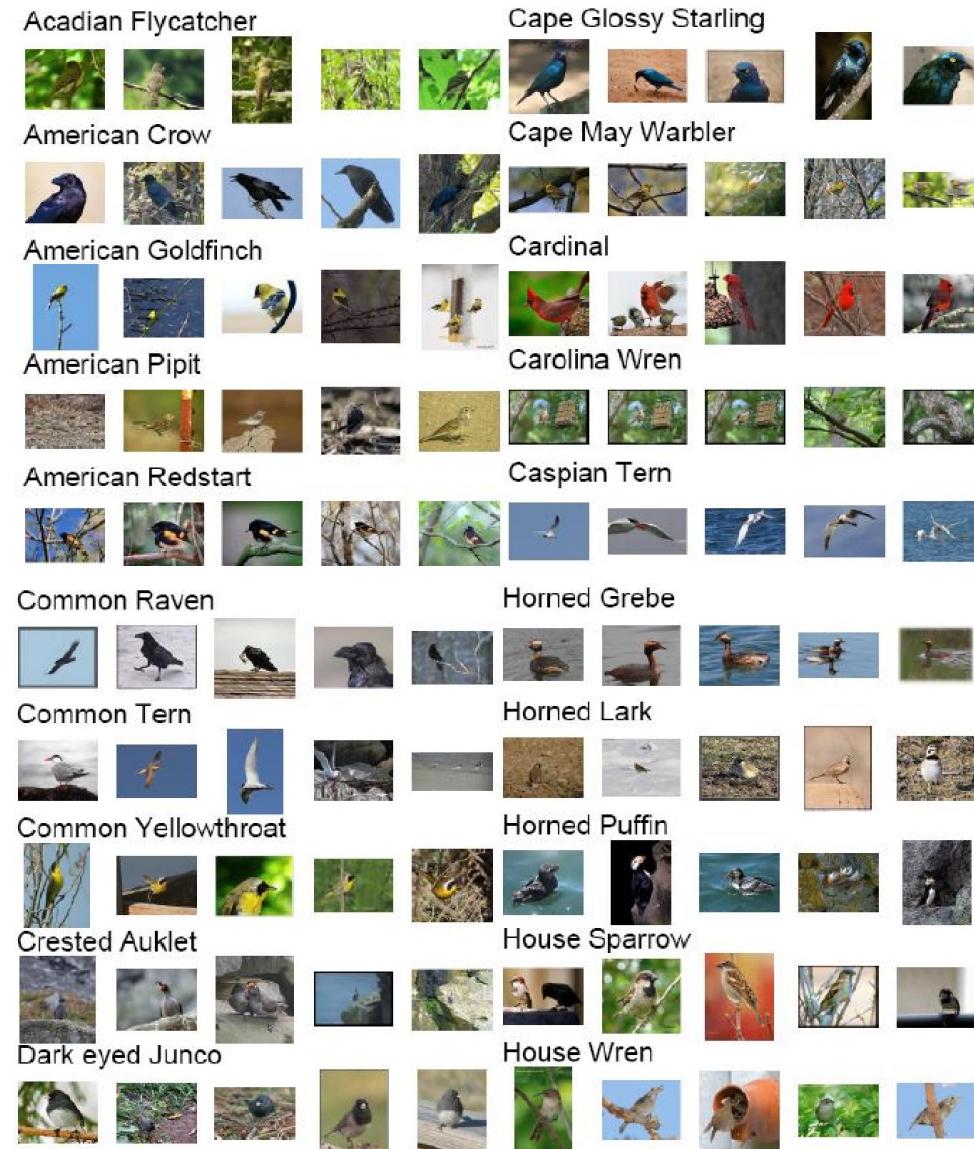
- Text-based image retrieval is useful
- If it's generalisable and highly domain-specific, even more so

- Excellent work by Reed et al. (2016)
- Let's revisit it using current SOTA
- Combine & fine-tune networks



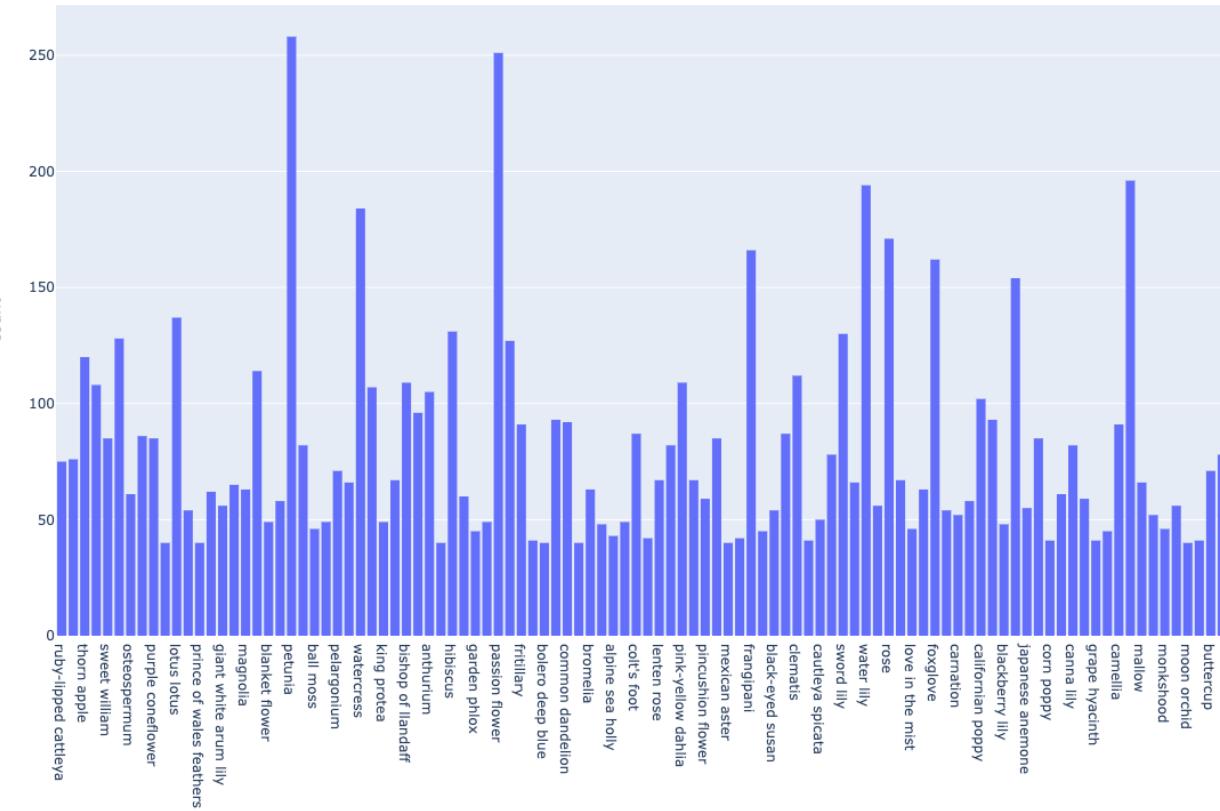
Datasets

- Caltech-UCSD Birds (+ 10 textual descriptions per image)
200 categories, 11,788 images
- Oxford-102 Flowers (+ 10 textual descriptions per image)
102 categories, 8,189 images
- ~200k (image, caption, species)
- Very fine-grained

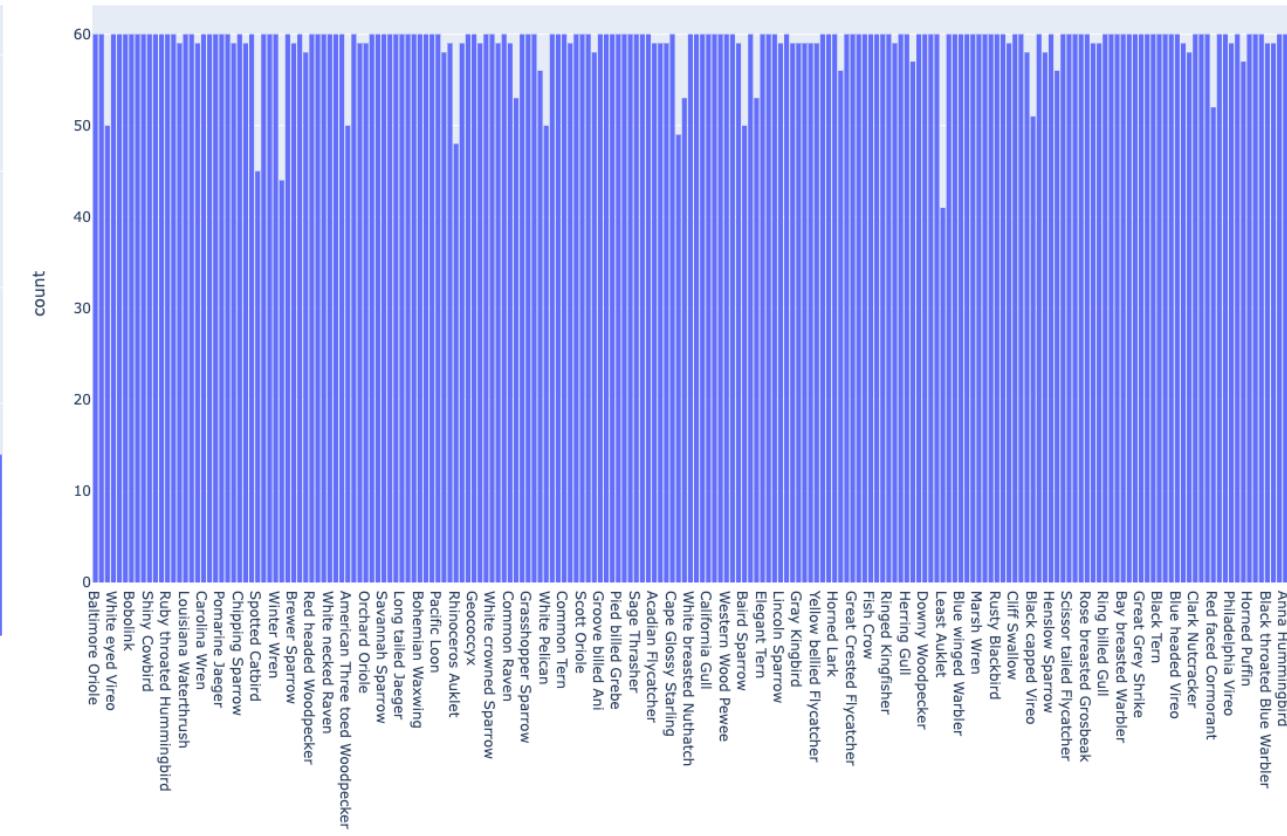


Datasets – class distribution

Flowers



Birds



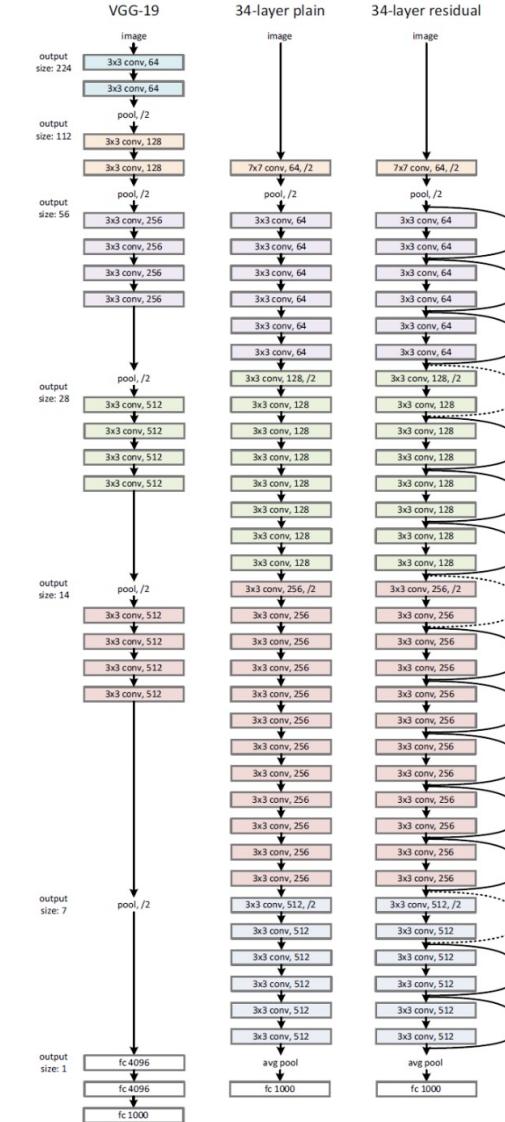
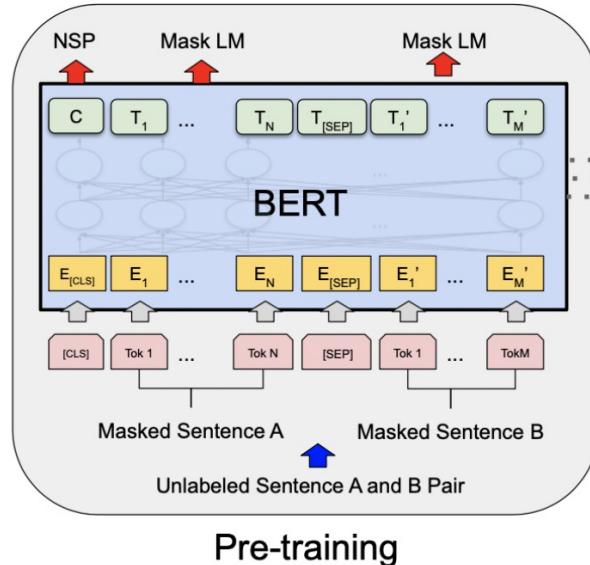
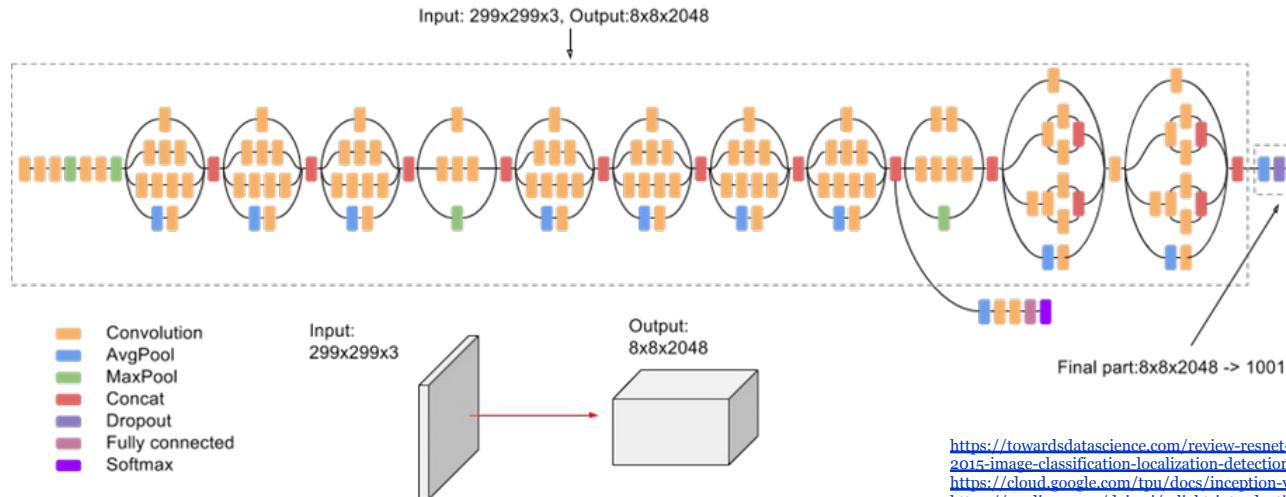
Background

- Text encoders:

- DistilBERT
- TF-IDF

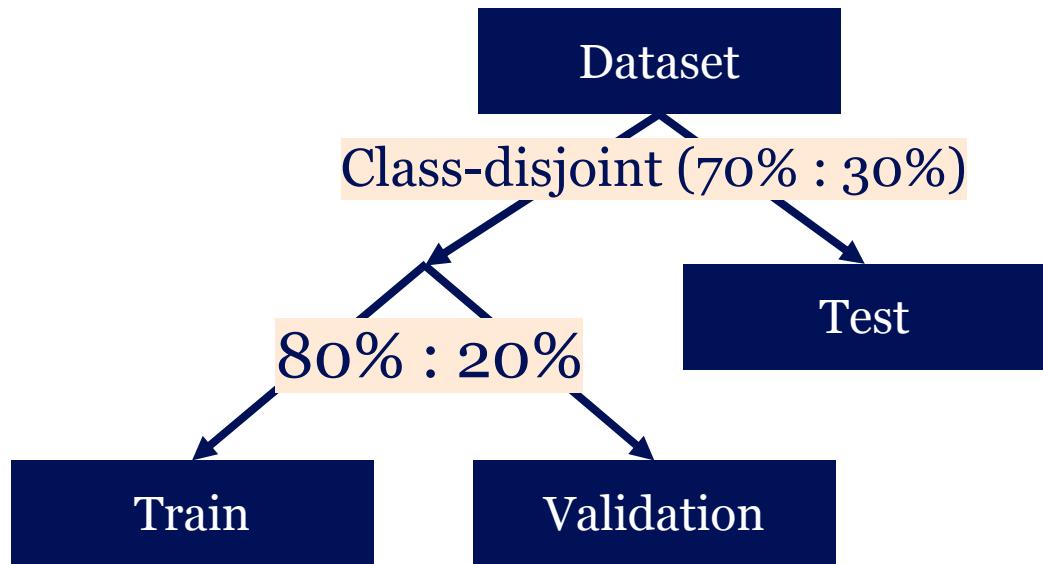
- Image classifiers:

- VGG
- ResNet
- Inception V3



<https://towardsdatascience.com/review-resnet-winner-of-ilsvrc-2015-image-classification-localization-detection-e20402bfa5d8>
<https://cloud.google.com/pml/docs/inception-v3-advanced>
<https://medium.com/dair-ai/a-light-introduction-to-bert-2da54f06b68c>

Methods



Methods – steps

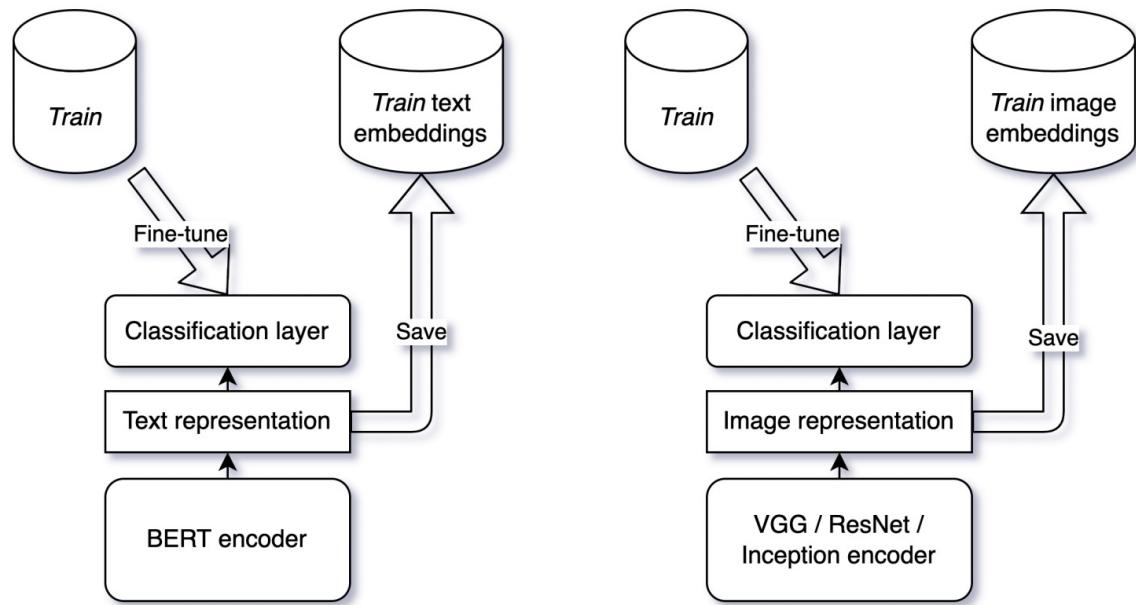


Fig. 3. The first training step is to separately fine-tune the networks via text-classification and image-classification. *Train* denotes a class-disjoint split of either *birds* or *flowers*.

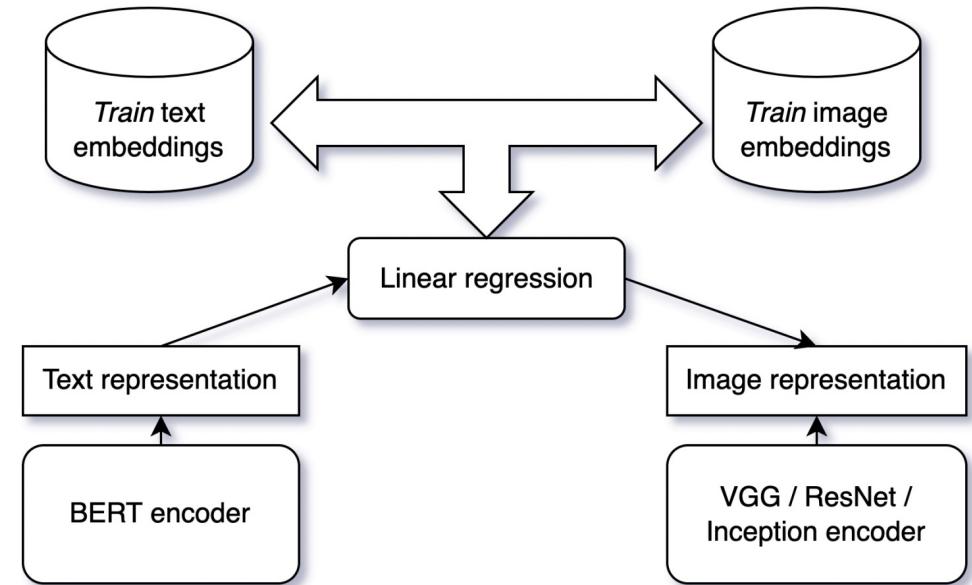


Fig. 4. For the second-stage of training, the last, classification layers are removed from both models. The exposed, raw text embeddings are mapped to the image-embeddings using a linear projection learned from the actual embeddings of the *train* split.

Methods – steps

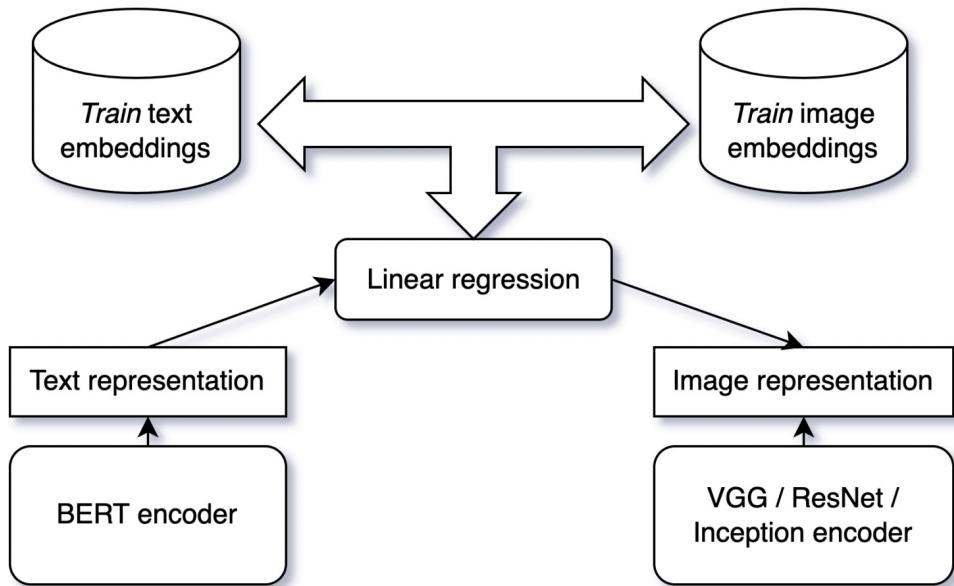


Fig. 4. For the second-stage of training, the last, classification layers are removed from both models. The exposed, raw text embeddings are mapped to the image-embeddings using a linear projection learned from the actual embeddings of the *train* split.

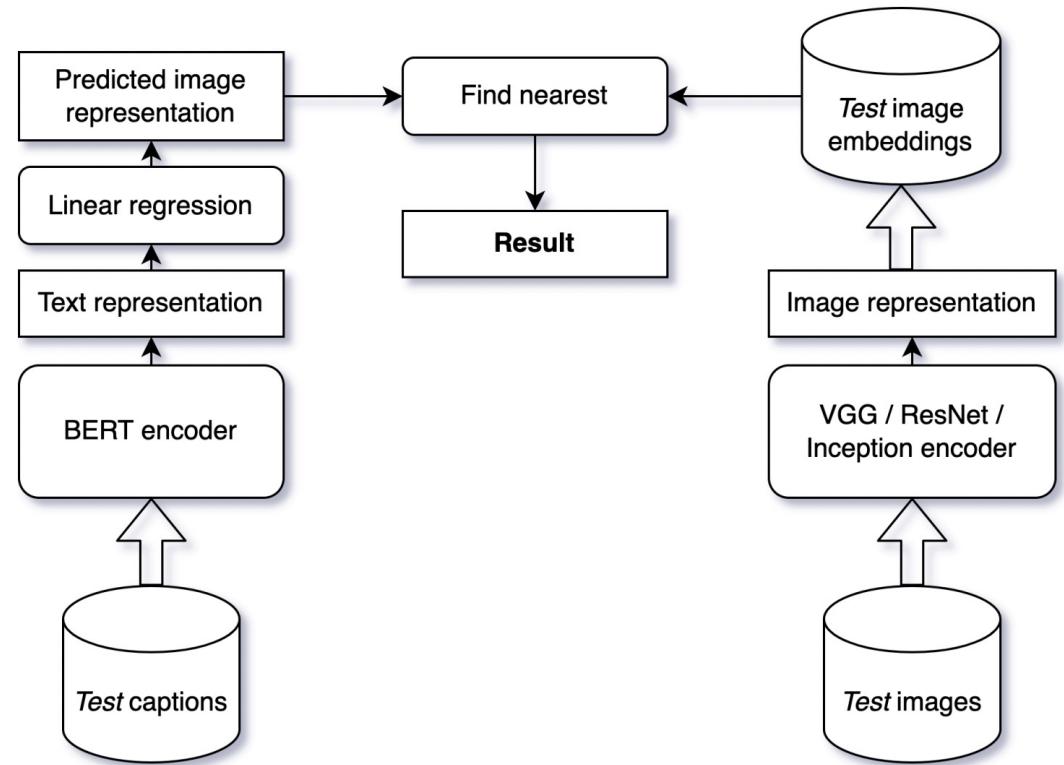


Fig. 5. First, each image from the *test* split is encoded using the fine-tuned and subsequently beheaded image-classifier. For each caption, an image embedding is predicted using the fine-tuned text encoder and a linear projection. Finally, the images with the most similar embedding to the predicted image embedding are returned.

Results – fine-tuning

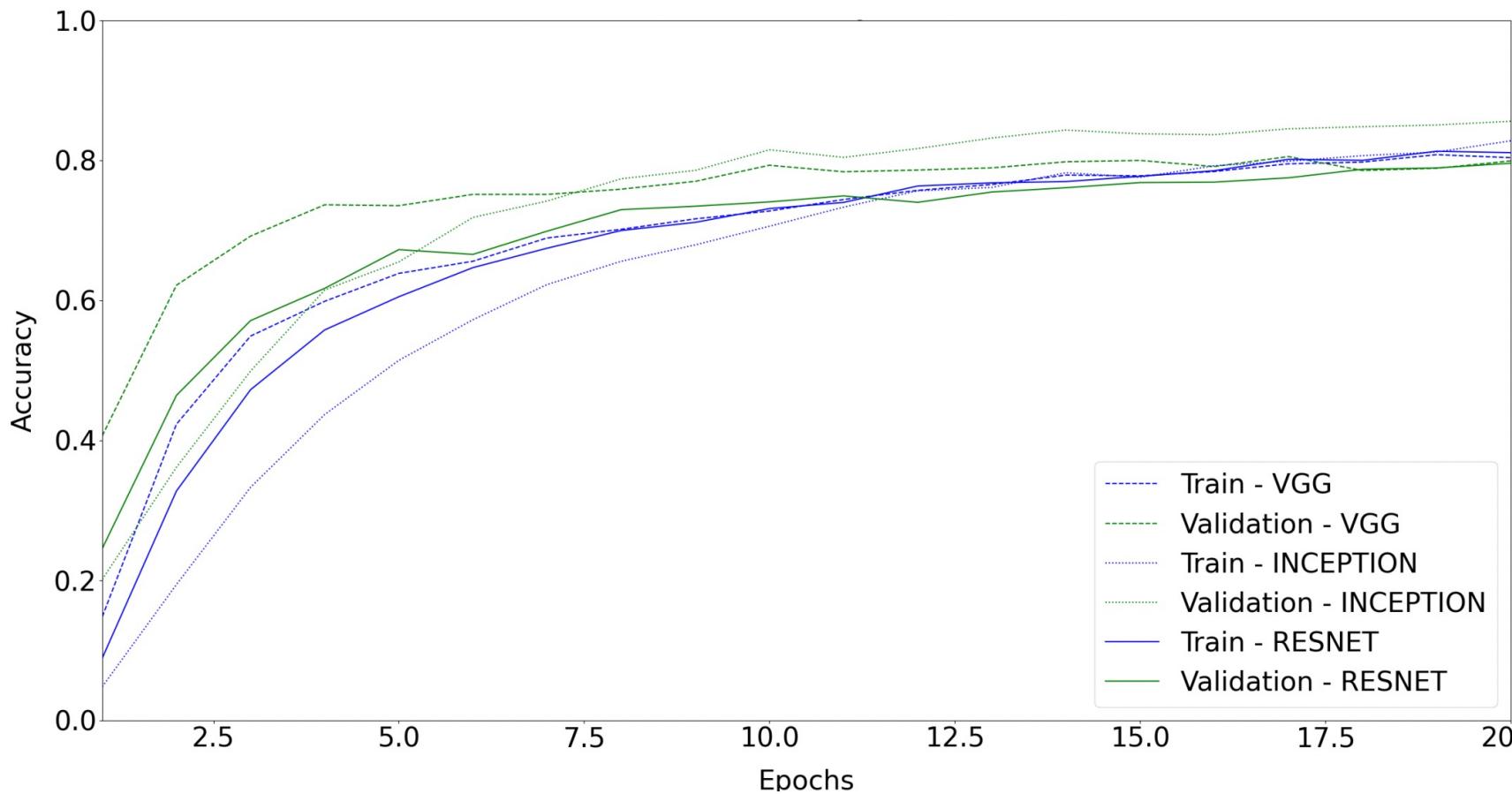


Fig. 6. The accuracy metrics over time of the image classifiers on the *birds* train and validation datasets.

Results – fine-tuning

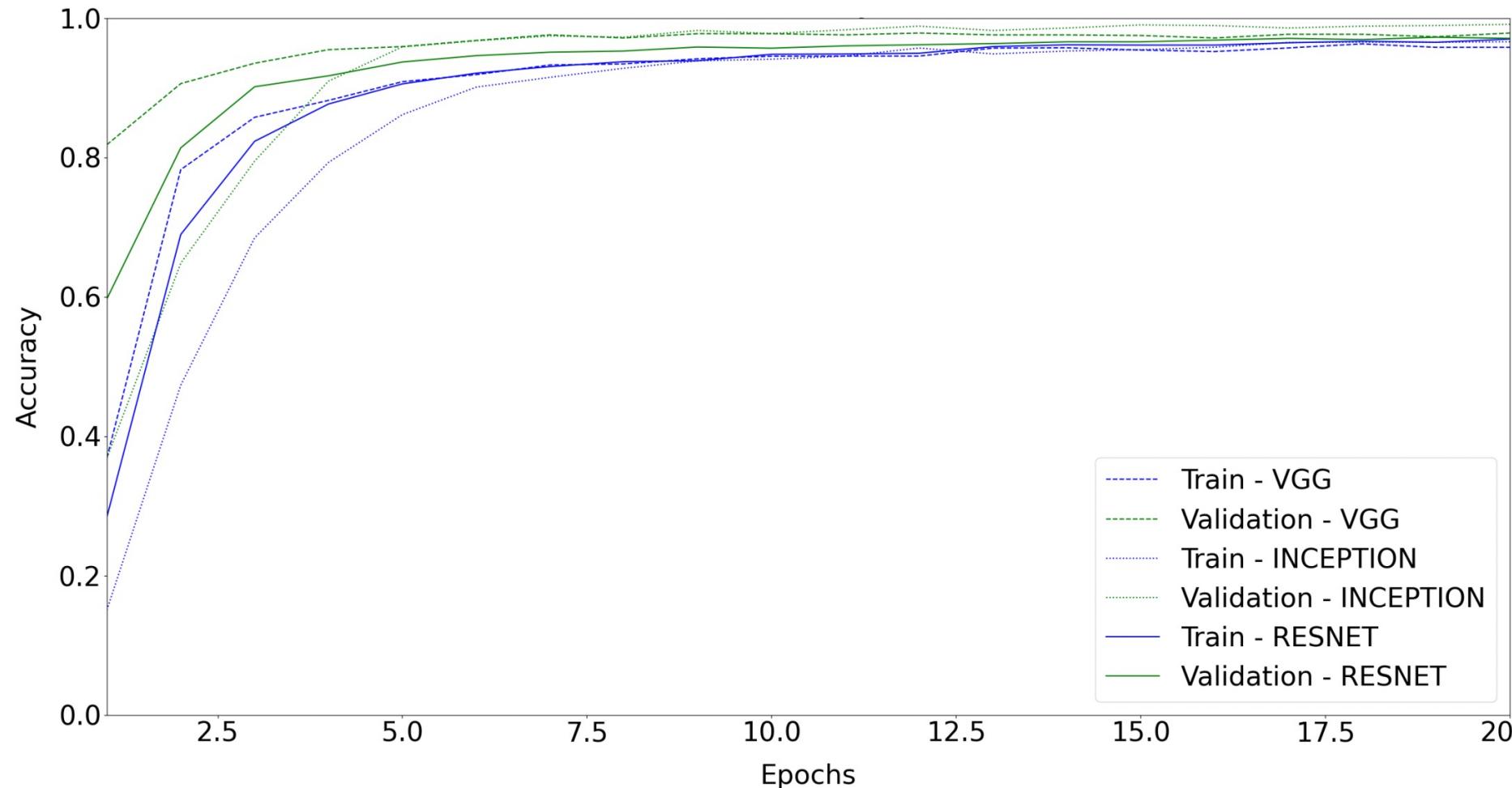


Fig. 7. The accuracy metrics over time of the image classifiers on the *flowers* train and validation datasets.

Results tables

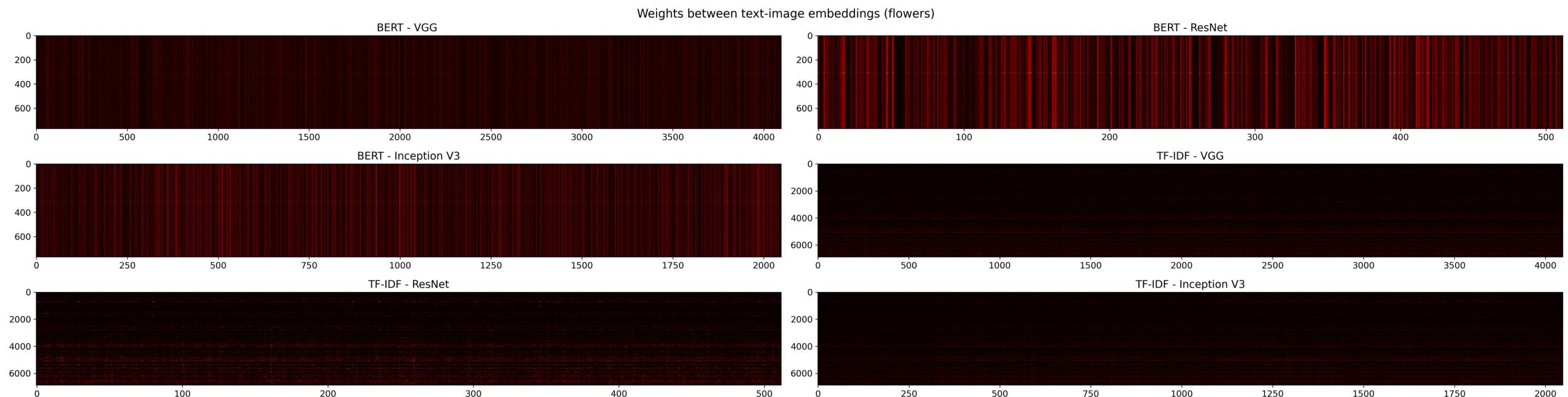
Table 1. Comparison of the 6 presented models on the task of zero-shot text-based image retrieval using the Caltech-UCSD *birds* 200 [10] dataset.

	Top-1 Accuracy	Precision	Recall	F1-score	Top-5 Accuracy
DistilBERT - VGG	0.19	0.25	0.19	0.19	0.43
DistilBERT - ResNet	0.11	0.17	0.11	0.10	0.32
DistilBERT - Inception v3	0.12	0.33	0.11	0.11	0.29
TF-IDF - VGG	0.11	0.28	0.11	0.11	0.33
TF-IDF - ResNet	0.10	0.33	0.10	0.09	0.24
TF-IDF - Inception v3	0.10	0.38	0.10	0.09	0.21

Table 2. Comparison of the 6 presented models on the task of zero-shot text-based image retrieval using the Oxford-102 *flowers* [11] dataset.

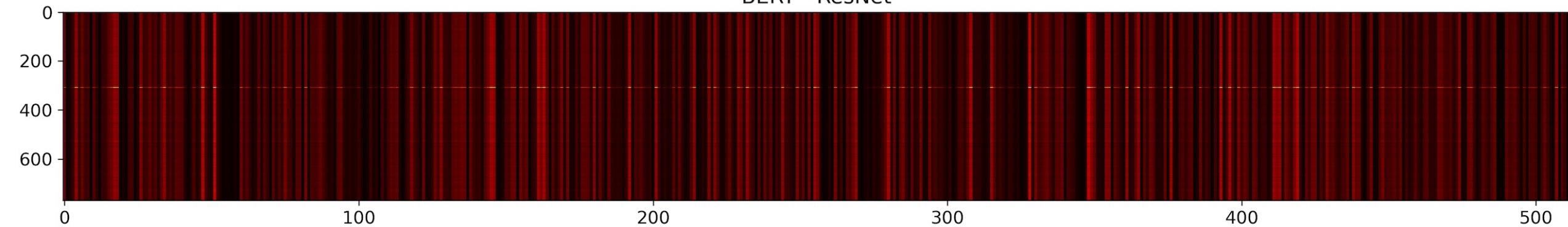
	Top-1 Accuracy	Precision	Recall	F1-score	Top-5 Accuracy
DistilBERT - VGG	0.37	0.30	0.27	0.23	0.57
DistilBERT - ResNet	0.36	0.31	0.26	0.24	0.53
DistilBERT - Inception v3	0.40	0.29	0.29	0.27	0.58
TF-IDF - VGG	0.20	0.28	0.16	0.13	0.46
TF-IDF - ResNet	0.18	0.28	0.13	0.12	0.33
TF-IDF - Inception v3	0.24	0.24	0.17	0.15	0.43

Results – weights

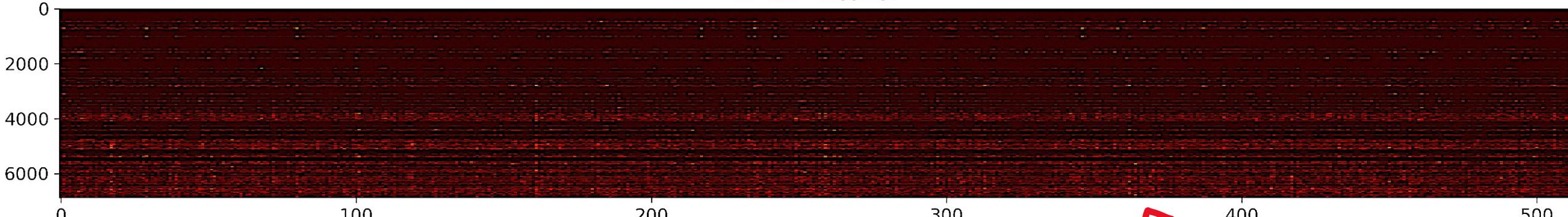


Results – weights

BERT - ResNet

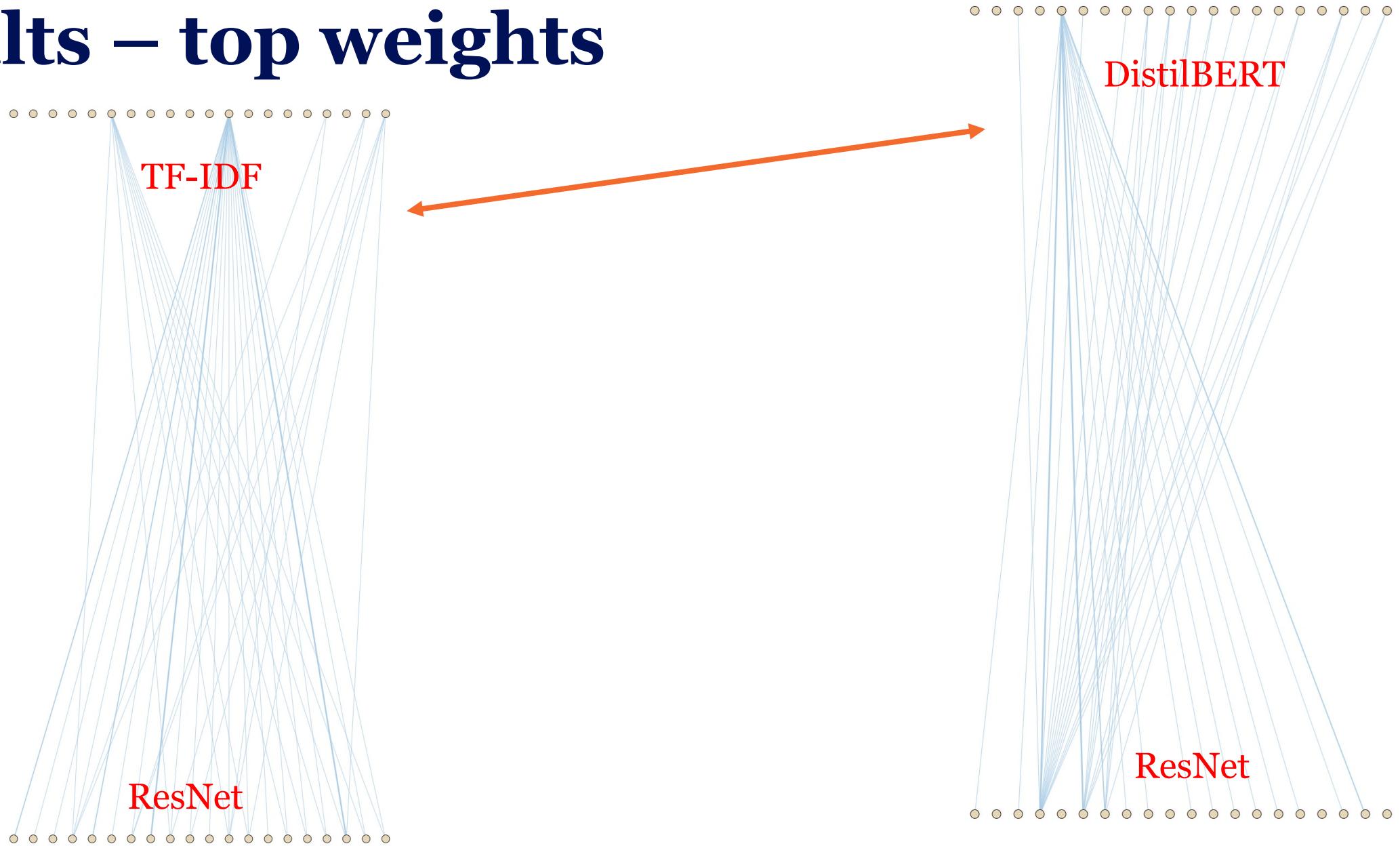


TF-IDF - ResNet



Colour-adjusted for better visibility;
see original in the report.

Results – top weights



Conclusion

- **6 architectures** were compared
- On highly domain-specific **fine-grained** classes, in a **zero-shot** setup
- Fine-tuned on different tasks then combined together
- **It worked**
- **There is similarity in the embedding spaces**
- TF-IDF << **DistilBERT**
- No absolute winner for image processing architectures
- **Worse** than Reed et al.'s **but more general**

Datasets

- Caltech-UCSD Birds (+ 10 textual descriptions per image)
200 categories, 11,788 images
- Oxford-102 Flowers (+ 10 textual descriptions per image)
102 categories, 8,189 images
- ~200k (image, caption, species)
- Very fine-grained



Discover the world at Leiden University

3

Methods – steps

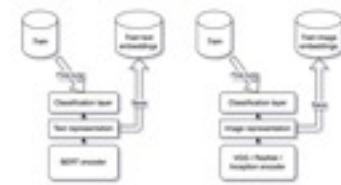


Fig. 3. The first training step is to separately fine-tune the networks via text-classification and image-classification. Train denotes a class-disjoint split of either birds or flowers.

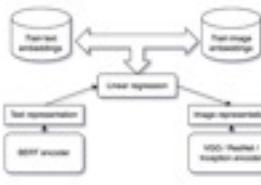


Fig. 4. For the second stage of training, the test classification layers are removed from both nets. The learned new test embeddings are mapped to the image embeddings using a linear projection learned from the actual embeddings of the main split.

Discover the world at Leiden University

Results – fine-tuning

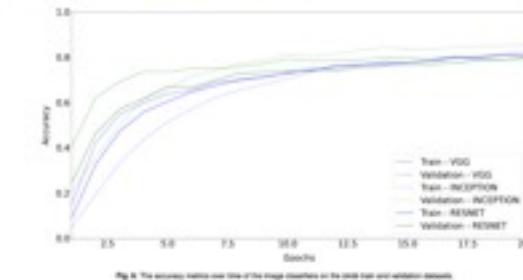


Fig. 5. The accuracy metric over time of the image classifiers on the Caltech train and validation datasets.

Discover the world at Leiden University

4

Results tables

Table 1. Comparison of the 6 presented models on the task of zero-shot text-based image retrieval using the Caltech-UCSD Birds [1] dataset.

	Top-1 Accuracy	Precision	Recall	F1-score	Top-5 Accuracy
DistilBERT - VGG	0.29	0.25	0.29	0.25	0.43
DistilBERT - ResNet	0.23	0.17	0.20	0.16	0.32
DistilBERT - Inception v3	0.17	0.10	0.11	0.09	0.29
TF-IDF - VGG	0.13	0.28	0.10	0.13	0.13
TF-IDF - ResNet	0.09	0.15	0.09	0.13	0.18
TF-IDF - Inception v3	0.10	0.10	0.09	0.11	0.21

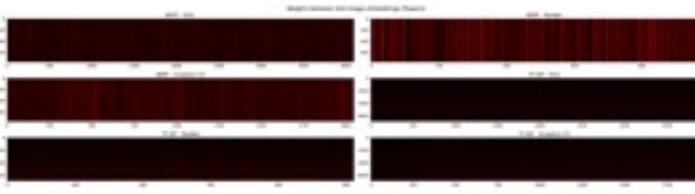
Table 2. Comparison of the 6 presented models on the task of zero-shot text-based image retrieval using the Oxford-102-Flowers [2] dataset.

	Top-1 Accuracy	Precision	Recall	F1-score	Top-5 Accuracy
DistilBERT - VGG	0.37	0.30	0.27	0.23	0.57
DistilBERT - ResNet	0.36	0.31	0.26	0.24	0.53
DistilBERT - Inception v3	0.40	0.29	0.29	0.27	0.58
TF-IDF - VGG	0.20	0.28	0.16	0.13	0.46
TF-IDF - ResNet	0.18	0.28	0.13	0.12	0.33
TF-IDF - Inception v3	0.24	0.24	0.17	0.15	0.43

Discover the world at Leiden University

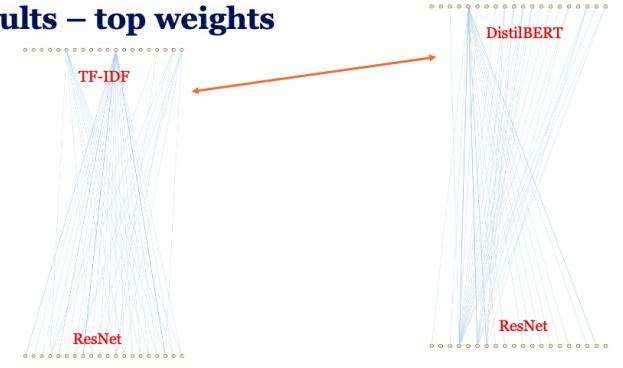
10

Results – weights



Discover the world at Leiden University

Results – top weights



Discover the world at Leiden University

14

<https://github.com/schmelczer/mir-final>



Universiteit
Leiden
The Netherlands

Discover the world at Leiden University

Thank you!