

# Few-shot Image Classification for Breast Cancer Detection

## An Advanced Approach in Medical Imaging

Qirun Dai, Jingzhi Sun,  
Pengyu Chen, Xiaowei Zeng

UC Davis & Fudan University

November 30, 2023

# Problem Statement

- **Background:** Breast cancer is a pressing concern worldwide, ranking as the second leading cause of cancer death in women. Early detection and identification of breast cancer can lead to timely treatment, effectively reducing the risk of further deterioration or death [Miller et al. 2022].
- **Motivation:**
  - Traditional medical image classification: Consumes a lot of manpower and time.
  - Pattern recognition and machine learning: Automates and streamlines medical image classification.
- **Objective:** Few-shot image classification for breast cancer detection.

# Major Challenges and Guidelines

- **Data sparsity:** How to generalize the model and prevent overfitting in a few-shot training setting? [Varoquaux and Cheplygina 2022]
  - Utilize a special breast image dataset that aligns with a genuine few-shot training scenario.
- **Rapid development of new vision models:** How to compare and evaluate the traditional machine learning models and cutting-edge deep learning models?
  - **Parametric supervised learning:** Naïve Bayes, Logistic Regression.
  - **Non-parametric supervised learning:** Support Vector Machine (SVM).
  - **Ensembling methods:** Decision Tree, Random Forest.
  - **Deep learning models:** Convolutional Neural Networks (CNN), Vision Transformers (ViT).

# Breast Ultrasound Images Dataset

- **Collection:** 2018 [Al-Dhabayani et al. 2020].
- **Observations:** 600 breast ultrasound images among women in ages between 25 and 75 years old.
- **Data size:** 780 images with an average image size of 500\*500 pixels.
- **Labels:** Three classes, which are **normal, benign, and malignant**.
- **Format:** PNG, including original images and masked images.
- **Randomly split:** 80% for training and the remaining 20% for testing.
- This dataset not only vividly simulates the real-world scenario of medical data sparsity, but also poses a great challenge to the efficient training of machine learning models.

# Masking in Image Classification

**Masking** is a technique in image processing that highlights specific areas, crucial in medical imaging for focusing on regions of interest. This approach offers several advantages:

- **Enhanced Focus:** Directs analysis to key regions, aiding in accurate diagnosis.
- **Accuracy:** Improves classification accuracy by focusing on relevant areas (Shape and Area Size).
- **Noise Reduction:** Reduces background noise and extraneous details.
- **Efficiency:** Simplifies image content for more targeted analysis.

# Illustration of Masked Images

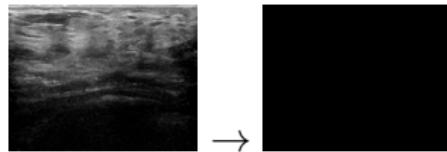


Figure: Normal and its Mask

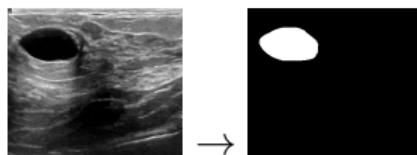


Figure: Benign and its Mask

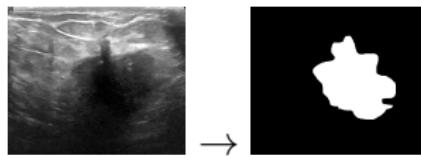


Figure: Malignant and its Mask

# Methods and Models

- ① Naïve Bayes (Parametric)
- ② Logistic Regression (Parametric)
- ③ Support Vector Machine (Non-Parametric)
- ④ Decision Tree and Random Forest (Ensemble)
- ⑤ Deep Learning Methods

# Naïve Bayes

- "**Naïve**": The strong assumption that the features are conditionally independent of one another given the class label, which enhances the algorithm's computational efficiency.
- **Performance in real-world scenarios:** While the independence assumption is often violated in practice, Naïve Bayes delivers competitive classification accuracy [Webb, Keogh, and Miikkulainen 2010].
- **Features:**
  - Computational efficiency.
  - Robustness in the face of noise.
  - Robustness in the face of missing values.

# Naïve Bayes in Medical Imaging

- Naïve Bayes is really **suitable** for medical image classification tasks.
- **Resilience against noise:**
  - It uses **all attributes** for predictions, regardless of the presence of noisy or irrelevant attributes.
  - It estimates the probability based on the **overall likelihood**.
- **Alignment with independence assumption:**
  - Due to the randomness of medical events [Zaw, Maneerat, and Win 2019].
  - Consequently, Naïve Bayes consistently delivered high predictive accuracy [Ramesh Kumar and Vijaya 2022].

# Naïve Bayes on Original Images

Images were resampled to a 64x64 resolution. We trained the Naïve Bayes classifier with default hyperparameters as the baseline.

| Confusion Matrix | Benign | Malignant | Normal |
|------------------|--------|-----------|--------|
| Benign           | 45     | 12        | 27     |
| Malignant        | 8      | 25        | 10     |
| Normal           | 12     | 6         | 11     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.69      | 0.54   | 0.60     | 84      |
| Malignant | 0.58      | 0.58   | 0.58     | 43      |
| Normal    | 0.23      | 0.38   | 0.29     | 29      |

Accuracy: 0.52    Total Support: 156

Table: Classification Report

# Naïve Bayes on Masked Images

We trained another Naïve Bayes classifier with default parameters on masked images. The following result is the baseline for later parameter tuning.

| Confusion Matrix | Benign | Malignant | Normal |
|------------------|--------|-----------|--------|
| Benign           | 79     | 8         | 0      |
| Malignant        | 30     | 16        | 0      |
| Normal           | 0      | 0         | 27     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.72      | 0.91   | 0.81     | 87      |
| Malignant | 0.67      | 0.35   | 0.46     | 46      |
| Normal    | 1.00      | 1.00   | 1.00     | 27      |

Accuracy: 0.76 Total Support: 160

Table: Classification Report

# Parameter Tuning for Naïve Bayes

`var_smoothing` is the stability calculation to smooth the curve and therefore account for more samples that are further away from the distribution mean. After the log-scale search and cross validation, we chose 0.0032.

| Confusion Matrix | Benign | Malignant | Normal |
|------------------|--------|-----------|--------|
| Benign           | 78     | 9         | 0      |
| Malignant        | 24     | 22        | 0      |
| Normal           | 0      | 0         | 27     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.76      | 0.90   | 0.83     | 87      |
| Malignant | 0.71      | 0.48   | 0.57     | 46      |
| Normal    | 1.00      | 1.00   | 1.00     | 27      |

Accuracy: 0.79 Total Support: 160

Table: Classification Report

# Findings in Naïve Bayes Learning Process

During the parameter tuning and evaluation of the Naïve Bayes model, we made several key observations:

- ① **Sampling Adequacy:** A 64x64 resampling was found to be sufficient for this problem. Increasing the sampling dimensions did not contribute to higher accuracy.
- ② **Computational Efficiency:** The training time is **linear** with both the number of training examples and the number of attributes, and the classification time is **linear** with the number of attributes and unaffected by the number of training examples.
- ③ **Improvement After Parameter Tuning:** Higher **recall** rate for malignant tumor and higher overall accuracy.
- ④ **Classification Challenges:** Models sometimes misclassify benign and malignant tumors, requiring more effective identification

# Methods and Models

- ① Naïve Bayes (Parametric)
- ② **Logistic Regression (Parametric)**
- ③ Support Vector Machine (Non-Parametric)
- ④ Decision Tree and Random Forest (Ensemble)
- ⑤ Deep Learning Methods

# Logistic Regression

- Logistic regression is a cornerstone in supervised learning method in classification with the following features:
  - **Applicability to Small Datasets:** Logistic regression exhibits proficiency with relatively small datasets, making it valuable in medical scenarios where acquiring labeled data brings privacy concerns and more cost.
  - **Linear Assumption:** Logistic regression assumes that the relationship between the log-odds of the variables should be approximately linear. While this assumption might be reasonable in certain cases, the complex nature of medical images introduces challenges to this linearity assumption.
- **Performance in medical image classification:** In a study by Dinesh and Kalyanasundaram 2022, logistic regression demonstrated superior performance compared to SVM, KNN, Decision Tree using the Wisconsin dataset.

# Logistic Regression on Original Images

Images were resampled to a 64x64 resolution. We trained the support vector machine with default hyperparameters and original images as the baseline.

| Confusion Matrix |           | Benign | Malignant | Normal |
|------------------|-----------|--------|-----------|--------|
|                  | Benign    | 71     | 8         | 5      |
|                  | Malignant | 17     | 24        | 2      |
|                  | Normal    | 15     | 2         | 12     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.69      | 0.85   | 0.76     | 84      |
| Malignant | 0.71      | 0.56   | 0.62     | 43      |
| Normal    | 0.63      | 0.41   | 0.50     | 29      |

Accuracy: 0.69 Total Support: 156

Table: Classification Report

# Logistic Regression on Masked Images

We trained another Logistic regression classifier with default parameters on masked images. The following result is the baseline for later parameter tuning.

| Confusion Matrix | Benign | Malignant | Normal |
|------------------|--------|-----------|--------|
| Benign           | 61     | 22        | 4      |
| Malignant        | 27     | 19        | 0      |
| Normal           | 0      | 0         | 27     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.69      | 0.70   | 0.70     | 87      |
| Malignant | 0.46      | 0.41   | 0.44     | 46      |
| Normal    | 0.87      | 1.00   | 0.93     | 27      |

Accuracy: 0.67 Total Support: 160

Table: Classification Report

# Parameter Tuning for Logistic Regression

Using grid search, we optimized the Logistic Regression classifier's hyperparameters.

**C:** 0.025, **penalty:** l2.

| Confusion Matrix | Benign | Malignant | Normal |
|------------------|--------|-----------|--------|
| Benign           | 72     | 14        | 1      |
| Malignant        | 23     | 23        | 0      |
| Normal           | 0      | 0         | 27     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.76      | 0.83   | 0.79     | 87      |
| Malignant | 0.62      | 0.50   | 0.55     | 46      |
| Normal    | 0.96      | 1.00   | 0.98     | 27      |

Accuracy: 0.76 Total Support: 160

Table: Classification Report

## 2-stage Logistic Regression

Introducing a mask image improved the classification accuracy between normal and tumors. Then train another Logistic Regression using original image to classify benign and malignant tumors

| Confusion Matrix | Benign | Malignant | Normal |
|------------------|--------|-----------|--------|
| Benign           | 83     | 8         | 1      |
| Malignant        | 17     | 22        | 0      |
| Normal           | 0      | 0         | 27     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.90      | 0.83   | 0.86     | 87      |
| Malignant | 0.77      | 0.73   | 0.75     | 46      |
| Normal    | 0.96      | 1.00   | 0.98     | 27      |

Accuracy: 0.84 Total Support: 160

Table: Classification Report

# Findings in Logistic Regression learning progress

During the parameter tuning and evaluation of Logistic Regression classifier, we made several key observations:

- ① **Sampling Adequacy:** A 64x64 resampling was found to be sufficient for this problem. Increasing the sampling dimensions did not contribute to higher accuracy.
- ② **Mask image:** Introducing a mask image significantly improved the classification accuracy between normal images and those with tumors, which led to a decrease in the classification accuracy between benign tumors and malignant tumors. So we introduce a two-stage Logistic Regression, which significantly improve the accuracy

# Methods and Models

- ① Naïve Bayes (Parametric)
- ② Logistic Regression (Parametric)
- ③ **Support Vector Machine (Non-Parametric)**
- ④ Decision Tree and Random Forest (Ensemble)
- ⑤ Deep Learning Methods

# Support Vector Machine

- **Theoretical Foundation:** Support Vector Machines (SVM) originated in the 1990s and, similar to an extension of neural networks, focuses on discovering an optimal hyperplane for classification.
- **Performance in medical images:** Chi, Feng, and Bruzzone 2008 emphasize SVM's efficiency in classifying small-sized training datasets, showcasing its ability to generalize well in scenarios with high-dimensional input spaces.
- **Key concepts:**
  - **Support Vector:** SVM relies on support vectors lying on class boundaries, which is useful in small datasets. This attribute proves especially beneficial in medical imaging.
  - **Kernel function:** SVM utilizes a kernel function to transform data into a higher-dimensional space, enabling the establishment of a decision plane to separate distinct classes.

# SVM on Original Images

Images were resampled to a 64x64 resolution. We trained the support vector machine with default hyperparameters and original images as the baseline.

| Confusion Matrix |           | Benign | Malignant | Normal |
|------------------|-----------|--------|-----------|--------|
|                  | Benign    | 78     | 6         | 0      |
|                  | Malignant | 21     | 22        | 0      |
|                  | Normal    | 22     | 2         | 5      |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.64      | 0.93   | 0.76     | 84      |
| Malignant | 0.73      | 0.51   | 0.60     | 43      |
| Normal    | 1.00      | 0.17   | 0.29     | 29      |

Accuracy: 0.67 Total Support: 156

Table: Classification Report

# SVM on Masked Images

We trained another support vector machine with default parameters on masked images. The following result is the baseline for later parameter tuning.

| Confusion Matrix | Benign | Malignant | Normal |
|------------------|--------|-----------|--------|
| Benign           | 77     | 7         | 3      |
| Malignant        | 16     | 30        | 0      |
| Normal           | 0      | 0         | 27     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.83      | 0.89   | 0.86     | 87      |
| Malignant | 0.81      | 0.65   | 0.72     | 46      |
| Normal    | 0.90      | 1.00   | 0.95     | 27      |

Accuracy: 0.84 Total Support: 160

Table: Classification Report

# Parameter Tuning for SVM

Using grid search, we optimized the support vector machine classifier's hyperparameters.

**C:** 22.5, **gamma:** 0.003, **kernel:** rbf.

| Confusion Matrix | Benign | Malignant | Normal |
|------------------|--------|-----------|--------|
| Benign           | 77     | 10        | 0      |
| Malignant        | 14     | 32        | 0      |
| Normal           | 0      | 0         | 27     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.85      | 0.89   | 0.87     | 87      |
| Malignant | 0.76      | 0.70   | 0.73     | 46      |
| Normal    | 1.00      | 1.00   | 1.00     | 27      |

Accuracy: 0.85 Total Support: 160

Table: Classification Report

# Findings in SVM learning progress

During the parameter tuning and evaluation of our support vector machine, we made several key observations:

- ① **Sampling Adequacy:** A 64x64 resampling was found to be sufficient for this problem.
- ② **Model performance:**
  - Improving accuracy from 0.67 to 0.84 through the inclusion of mask data is a significant enhancement, indicating that the additional information captured by the mask features is valuable for classification task.
  - Limited improvement of 0.01 through parameter tuning suggests that the model might already be performing close to its optimal capacity. Other factors are needed to improve performance.
  - Low False Negative Rate is critically important in medical image classification, especially when dealing with conditions like cancer.
- ③ **Classification Challenges:** SVM classifier sometimes misclassifies benign and malignant tumors, requiring more effective features.

# Methods and Models

- ① Naïve Bayes (Parametric)
- ② Logistic Regression (Parametric)
- ③ Support Vector Machine (Non-Parametric)
- ④ **Decision Tree and Random Forest (Ensemble)**
- ⑤ Deep Learning Methods

# Decision Tree and Random Forest

- **Decision Trees in Medical Imaging:** Kaganov, Ades, and Fraser 2018 highlighting the effectiveness of decision tree models in MRI signal intensity classification.
- **Advancement with Random Forests:** Criminisi et al. 2010 Describing the use of random forests for automatic detection and localization in three-dimensional CT scans. Increase accuracy and robust.
- **Study Aim:** Exploring the application of Random Forests in few-shot image classification for breast cancer detection - a three-class problem: normal, benign, or malignant.

## Initial Method with Raw Images

In this initial approach, images were resampled to a 64x64 resolution, training the Random Forest classifier with default hyperparameters. The obtained results were:

| Confusion Matrix |           | Benign | Malignant | Normal |
|------------------|-----------|--------|-----------|--------|
|                  | Benign    | 82     | 2         | 0      |
|                  | Malignant | 24     | 19        | 0      |
|                  | Normal    | 20     | 1         | 8      |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.65      | 0.98   | 0.78     | 84      |
| Malignant | 0.86      | 0.44   | 0.58     | 43      |
| Normal    | 1.00      | 0.28   | 0.43     | 29      |

Accuracy: 0.70 Total Support: 156

Table: Classification Report

# Results with Masked Images

After applying a 64x64 resampling and using default hyperparameters, the Random Forest classifier was trained on masked images. The results were as follows:

| Confusion Matrix |           | Benign | Malignant | Normal |
|------------------|-----------|--------|-----------|--------|
|                  | Benign    | 77     | 8         | 2      |
|                  | Malignant | 22     | 24        | 0      |
|                  | Normal    | 0      | 0         | 27     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.78      | 0.89   | 0.83     | 87      |
| Malignant | 0.75      | 0.52   | 0.62     | 46      |
| Normal    | 0.93      | 1.00   | 0.96     | 27      |

Accuracy: 0.80 Total Support: 160

Table: Classification Report

# Parameter Tuning with Grid Search

Using grid search, we optimized the Random Forest classifier's hyperparameters.

**Max Depth: 9, Min Sample Split: 4, Num of Estimators: 70.**

| Confusion Matrix | Benign | Malignant | Normal |
|------------------|--------|-----------|--------|
| Benign           | 81     | 3         | 3      |
| Malignant        | 17     | 29        | 0      |
| Normal           | 0      | 0         | 27     |

Table: Confusion Matrix

| Class     | Precision | Recall | F1-score | Support |
|-----------|-----------|--------|----------|---------|
| Benign    | 0.83      | 0.93   | 0.88     | 87      |
| Malignant | 0.91      | 0.63   | 0.74     | 46      |
| Normal    | 0.90      | 1.00   | 0.95     | 27      |

Accuracy: 0.86 Total Support: 160

Table: Classification Report

# Findings in Random Forest Learning Process

During the parameter tuning and evaluation of our Random Forest model, we made several key observations:

- ① **Sampling Adequacy:** A 64x64 resampling was found to be sufficient for this problem. Increasing the sampling dimensions did not contribute to higher accuracy.
- ② **Algorithm Stability:** The model showed good stability across different random seeds, indicating low overfitting.
- ③ **Classification Challenges:** Our masked-image model excelled in distinguishing diseased from non-diseased cases, likely due to non-diseased masks being entirely black while diseased ones show increased white values. However, it faced challenges differentiating between benign and malignant cases.

These findings provide insights into the strengths and limitations of our approach, guiding future improvements.

# Summary for Machine Learning Methods

The accuracy of the four machine learning methods are shown below.

| Accuracy        | Naïve Bayes | Logistic    | SVM         | Random Forest |
|-----------------|-------------|-------------|-------------|---------------|
| Original        | 0.52        | <b>0.69</b> | 0.67        | <b>0.70</b>   |
| Masked          | 0.76        | <b>0.67</b> | <b>0.84</b> | 0.80          |
| Masked + Tuning | 0.79        | 0.76        | 0.85        | <b>0.86</b>   |

Table: Summary of the Results

- Except the Logistic model, all the methods have higher accuracy when training with masked data than with original data.
- Random Forest and SVM performs better on breast cancer detection.
- The Logistic model has better performance in distinguishing between benign and malignant tumors when training with original data. While training with mask data, all models have better performance in distinguishing whether patient the has a cancer or not.  
( Normal v.s. Benign & Malignant).

# Methods and Models

- ① Naïve Bayes (Parametric)
- ② Logistic Regression (Parametric)
- ③ Support Vector Machine (Non-Parametric)
- ④ Decision Tree and Random Forest (Ensemble)
- ⑤ Deep Learning Methods

# **Deep Learning Methods**

## **Using Transfer Learning for Medical Image CLS**

# **Part 1:**

## **An Intuitive Introduction to Transfer Learning**

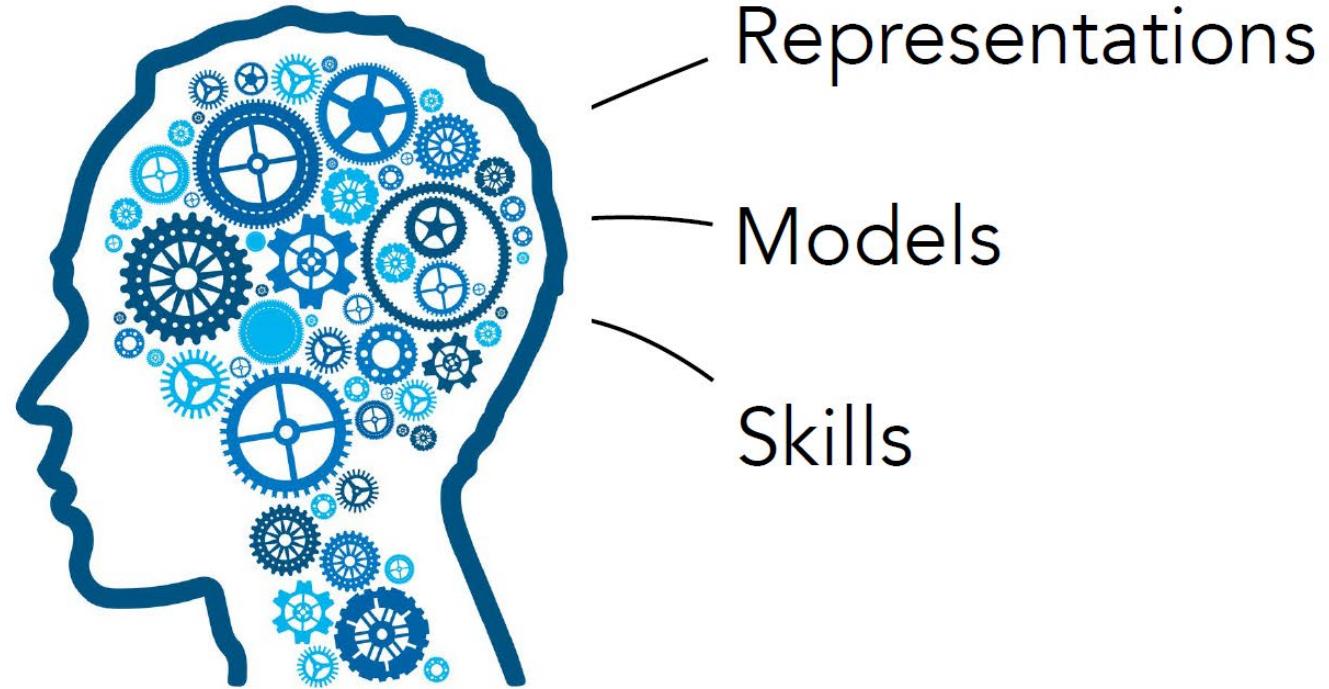
## **1.1 What is Transfer Learning?**

## **1.2 Why Transfer Learning?**

# "Deep learning"



# Human learning



How can we give deep nets prior knowledge?

# How can we give deep nets prior knowledge?

Just “pretrain” them on prior tasks!



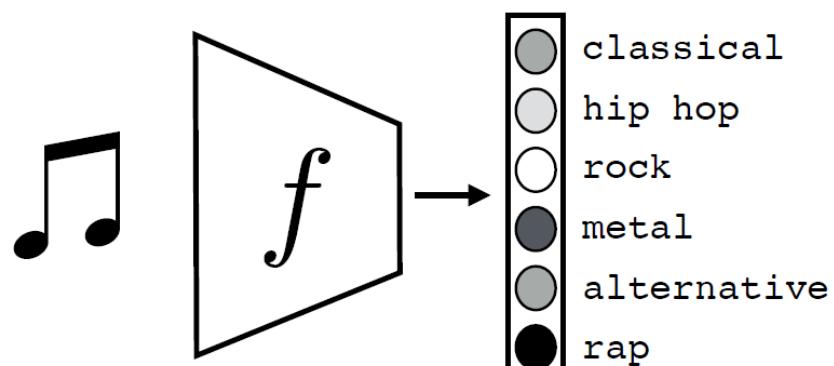
Then transfer knowledge from the pretrained nets to solve novel tasks.

One view: the **point** of deep learning is to enable problem solving with little data

## 1.3 How to do Transfer Learning?

## Pretraining

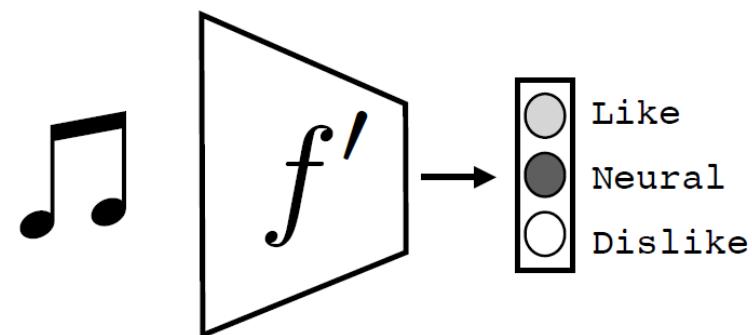
Genre recognition



*A lot of data*

## Adapting

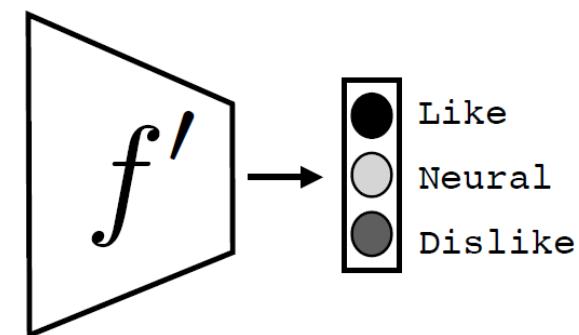
Preference prediction



*A little data*

## Testing

Preference prediction



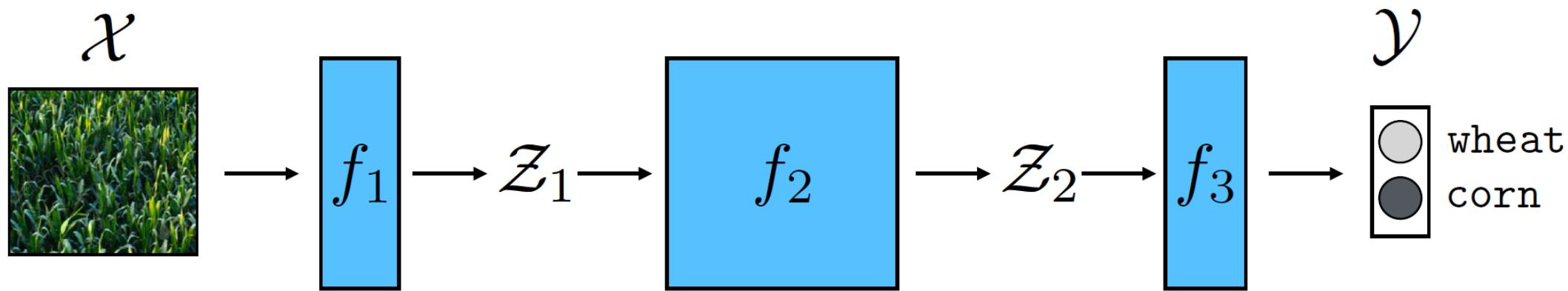
# Finetuning

- Pretrain a network on task A, resulting in parameters **W** and **b**
- Initialize a second network with some or all of **W** and **b**
- Train the second network on task B, resulting in parameters **W'** and **b'**

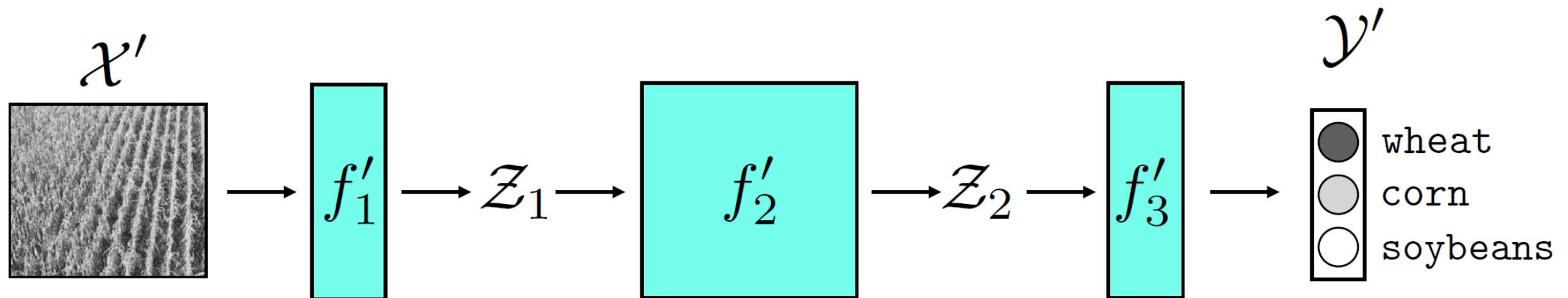
The “learned representation” is the encoder (its weights and biases) so that’s what we transfer

# What if the input/output dimensions don't match?

## Pretraining



## Finetuning



# **Part 2:**

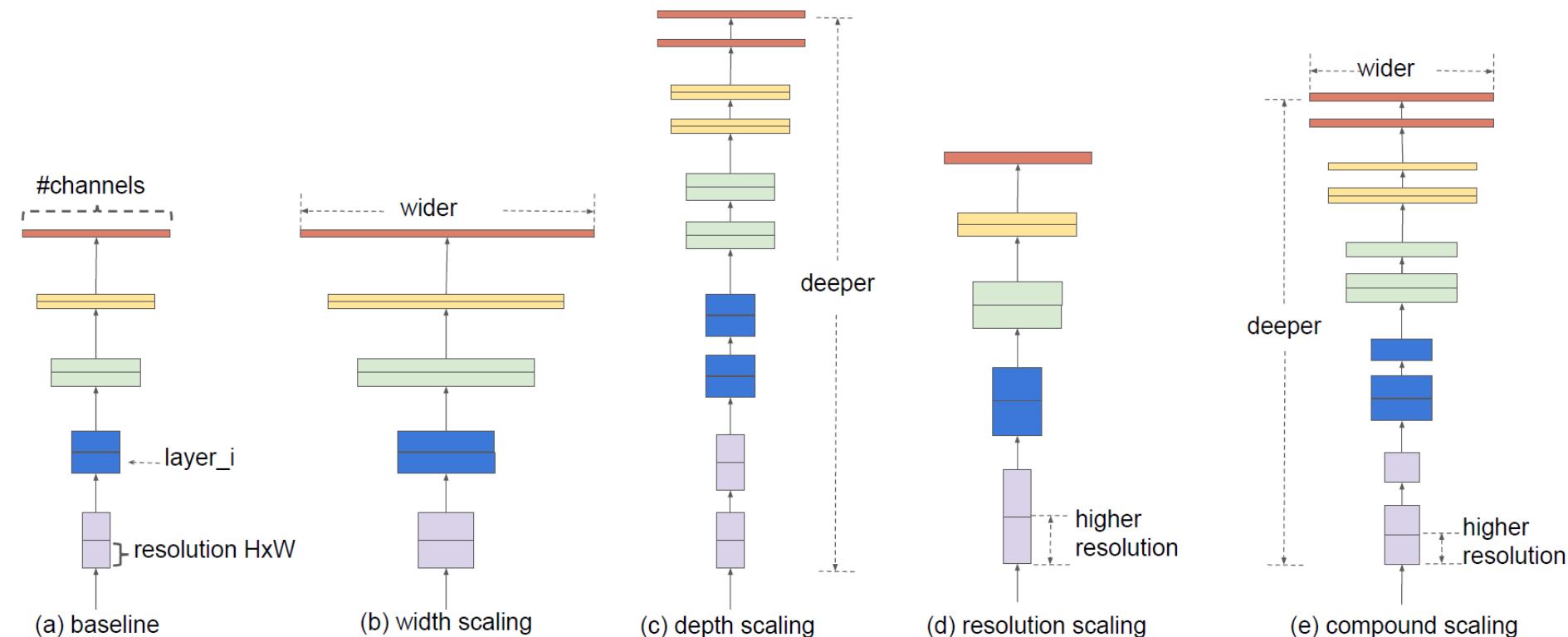
## **Results and Analysis of Transfer Learning in Breast Cancer Ultrasound Image Classification**

## **2.1 Base Model Architecture for Transfer Learning**

# 1) Convolutional Neural Networks (CNNs)

- EfficientNet (<https://arxiv.org/abs/1905.11946>; ICML 2019)
- A SOTA CNN model for Transfer Learning

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks



**Figure 2. Model Scaling.** (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is our proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio.

## 2) Vision Transformers (ViTs)

- Swin Transformer (<https://arxiv.org/abs/2103.14030> ; CVPR 2021)
- A SOTA ViT model for Transfer Learning

### Vision Transformer (ViT)

In practice: take 224x224 input image,  
divide into 14x14 grid of 16x16 pixel  
patches (or 16x16 grid of 14x14 patches)

Output vectors

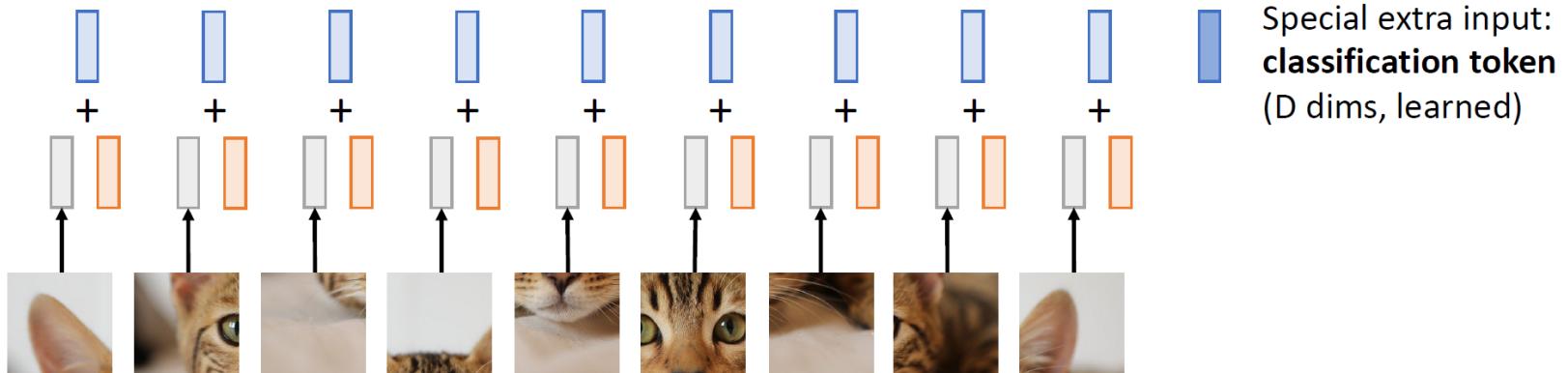


With 48 layers, 16 heads per  
layer, all attention matrices  
take 112 MB (or 192MB)

Linear projection  
to C-dim vector  
of predicted  
class scores

Exact same as  
NLP Transformer!

Add positional  
embedding: learned D-  
dim vector per position



---

**(a) Regular ImageNet-1K trained models**

| method           | image size       | #param. | FLOPs  | throughput (image / s) | ImageNet top-1 acc. |
|------------------|------------------|---------|--------|------------------------|---------------------|
| RegNetY-4G [48]  | 224 <sup>2</sup> | 21M     | 4.0G   | 1156.7                 | 80.0                |
| RegNetY-8G [48]  | 224 <sup>2</sup> | 39M     | 8.0G   | 591.6                  | 81.7                |
| RegNetY-16G [48] | 224 <sup>2</sup> | 84M     | 16.0G  | 334.7                  | 82.9                |
| EffNet-B3 [58]   | 300 <sup>2</sup> | 12M     | 1.8G   | 732.1                  | 81.6                |
| EffNet-B4 [58]   | 380 <sup>2</sup> | 19M     | 4.2G   | 349.4                  | 82.9                |
| EffNet-B5 [58]   | 456 <sup>2</sup> | 30M     | 9.9G   | 169.1                  | 83.6                |
| EffNet-B6 [58]   | 528 <sup>2</sup> | 43M     | 19.0G  | 96.9                   | 84.0                |
| EffNet-B7 [58]   | 600 <sup>2</sup> | 66M     | 37.0G  | 55.1                   | 84.3                |
| ViT-B/16 [20]    | 384 <sup>2</sup> | 86M     | 55.4G  | 85.9                   | 77.9                |
| ViT-L/16 [20]    | 384 <sup>2</sup> | 307M    | 190.7G | 27.3                   | 76.5                |
| DeiT-S [63]      | 224 <sup>2</sup> | 22M     | 4.6G   | 940.4                  | 79.8                |
| DeiT-B [63]      | 224 <sup>2</sup> | 86M     | 17.5G  | 292.3                  | 81.8                |
| DeiT-B [63]      | 384 <sup>2</sup> | 86M     | 55.4G  | 85.9                   | 83.1                |
| Swin-T           | 224 <sup>2</sup> | 29M     | 4.5G   | 755.2                  | 81.3                |
| Swin-S           | 224 <sup>2</sup> | 50M     | 8.7G   | 436.9                  | 83.0                |
| Swin-B           | 224 <sup>2</sup> | 88M     | 15.4G  | 278.1                  | 83.5                |
| Swin-B           | 384 <sup>2</sup> | 88M     | 47.0G  | 84.7                   | 84.5                |

---

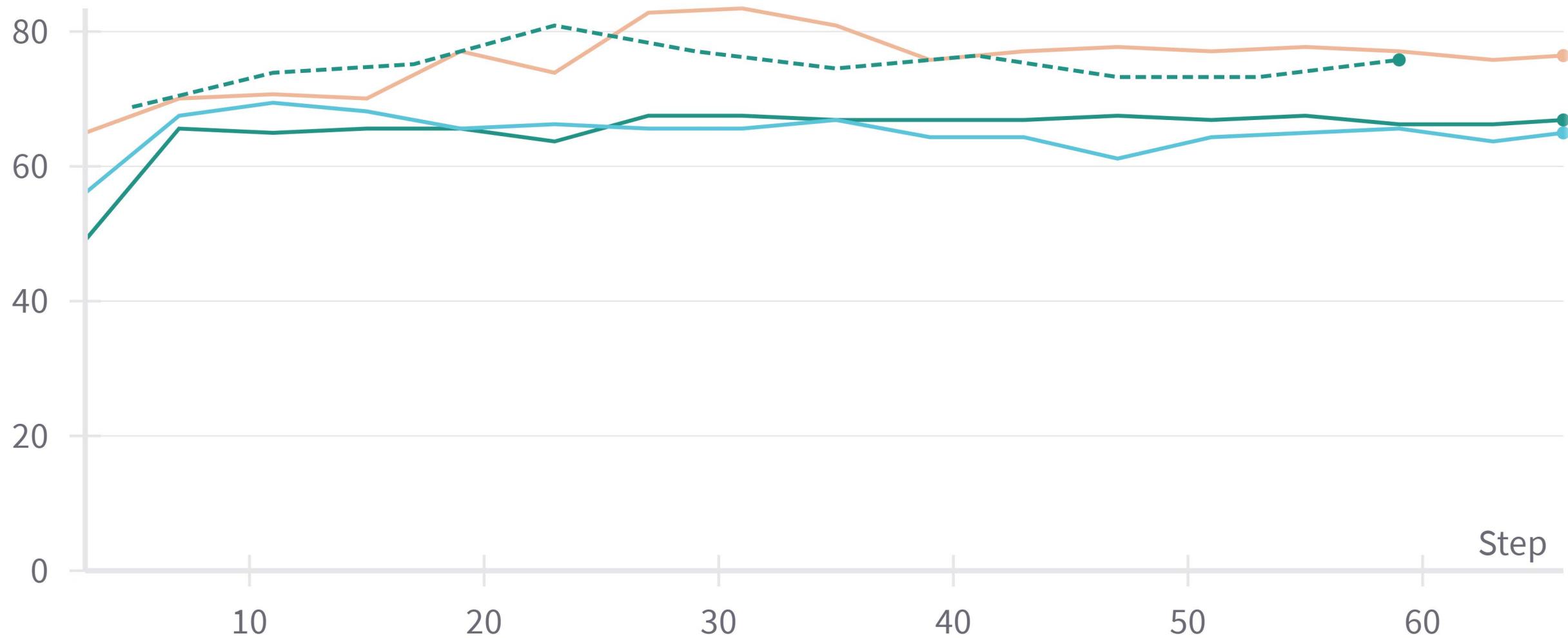
## **2.2 Final Performance Results**

# Best Performance Results (Only Raw Images; No Masking)

|              | Image Resolution | Data Augmentation | #Parameters | FLOPs | Top-1 Accuracy |
|--------------|------------------|-------------------|-------------|-------|----------------|
| Efficient_b5 | 224              | √                 | 30M         | 9.9G  | 69.4%          |
| Efficient_b7 |                  |                   | 66M         | 37.0G | 67.5%          |
| Swin_Tiny    |                  |                   | 29M         | 4.5G  | 80.9%          |
| Swin_Base    |                  |                   | 88M         | 15.4G | 83.4%          |

## dataset\_0\_test\_acc1

— swin-base-224-cosine lr=1e-4-num\_repeats=5    — efficientnetb7-224-cosine lr=1e-4-num\_repeats=5  
— efficientnetb5-224-cosine lr=1e-4-num\_repeats=5    -- swin-tiny-224-const lr=1e-4



## **2.3 Ablation Studies and Analysis**

# Ablation 1: Model Size (of the same architecture) & Accuracy

|                     | #Parameters | Top-1 Accuracy |
|---------------------|-------------|----------------|
| <b>Efficient_b5</b> | 30M         | <b>69.4</b>    |
| <b>Efficient_b7</b> | 66M         | 67.5           |

## Counterfactual:

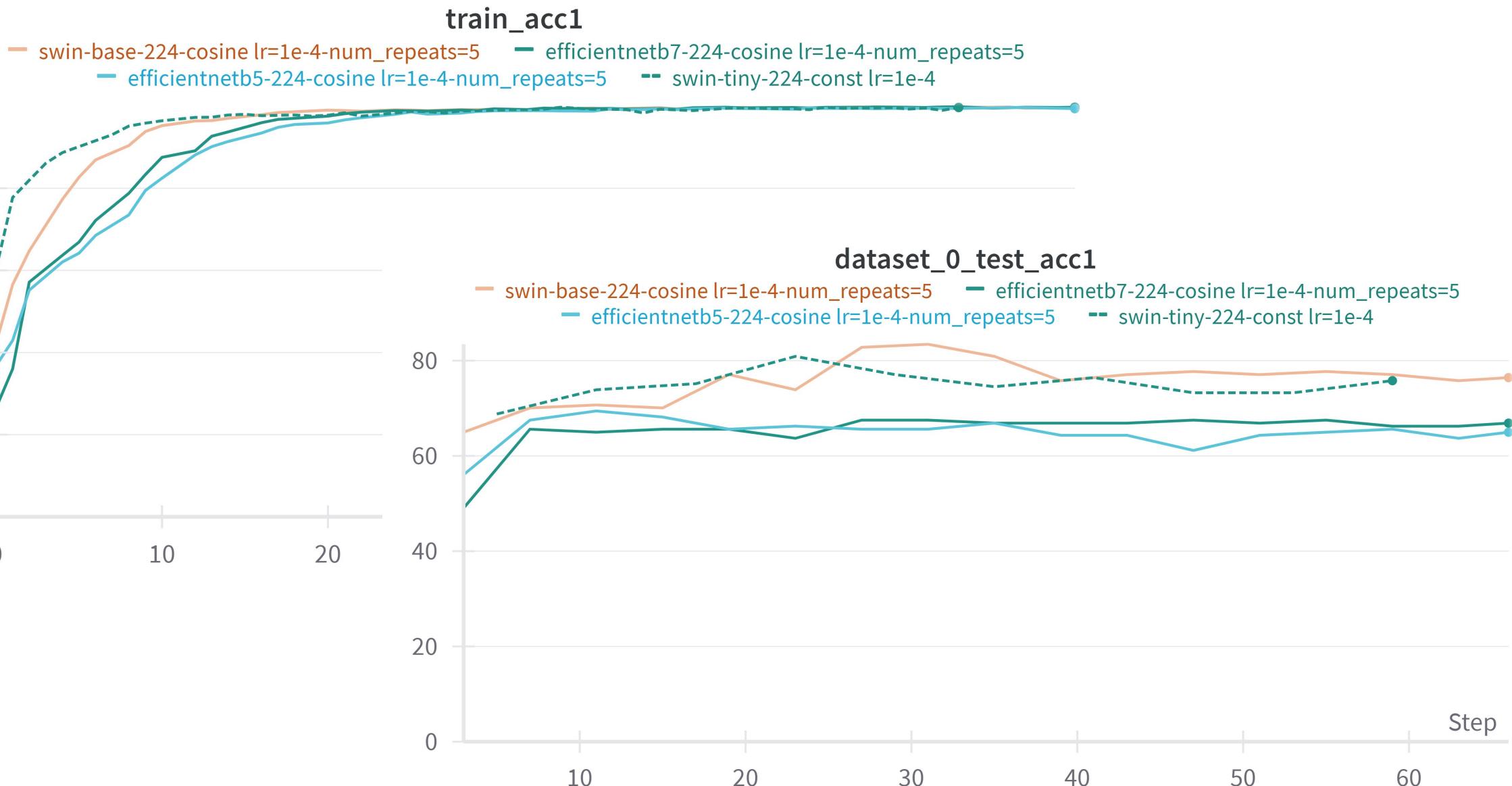
Efficient\_b5 is only half the size of Efficient\_b7, but still achieves better accuracy;

|                  | #Parameters | Top-1 Accuracy |
|------------------|-------------|----------------|
| <b>Swin_Tiny</b> | 29M         | 80.9           |
| <b>Swin_Base</b> | <b>88M</b>  | <b>83.4</b>    |

## On contrast:

Swin Base is larger than Swin Tiny, and the accuracy correspondingly scales.

# Explanation 1: Different Degrees of Overfitting



**But Another Question:**  
Why is ViT's Degree of Overfitting better than CNN?  
*(Will show the answer in Ablation 4)*

## Ablation 2: Image Resolution & Accuracy

|                     | <b>Image Resolution</b> | <b>Top-1 Accuracy</b> |
|---------------------|-------------------------|-----------------------|
| <b>Efficient_b5</b> | <b>224</b>              | <b>69.4</b>           |
| <b>Efficient_b5</b> | 456                     | 56.0                  |

|                  | <b>Image Resolution</b> | <b>Top-1 Accuracy</b> |
|------------------|-------------------------|-----------------------|
| <b>Swin_Base</b> | <b>224</b>              | <b>83.4</b>           |
| <b>Swin_Base</b> | 384                     | 81.5                  |

### Counterfactual:

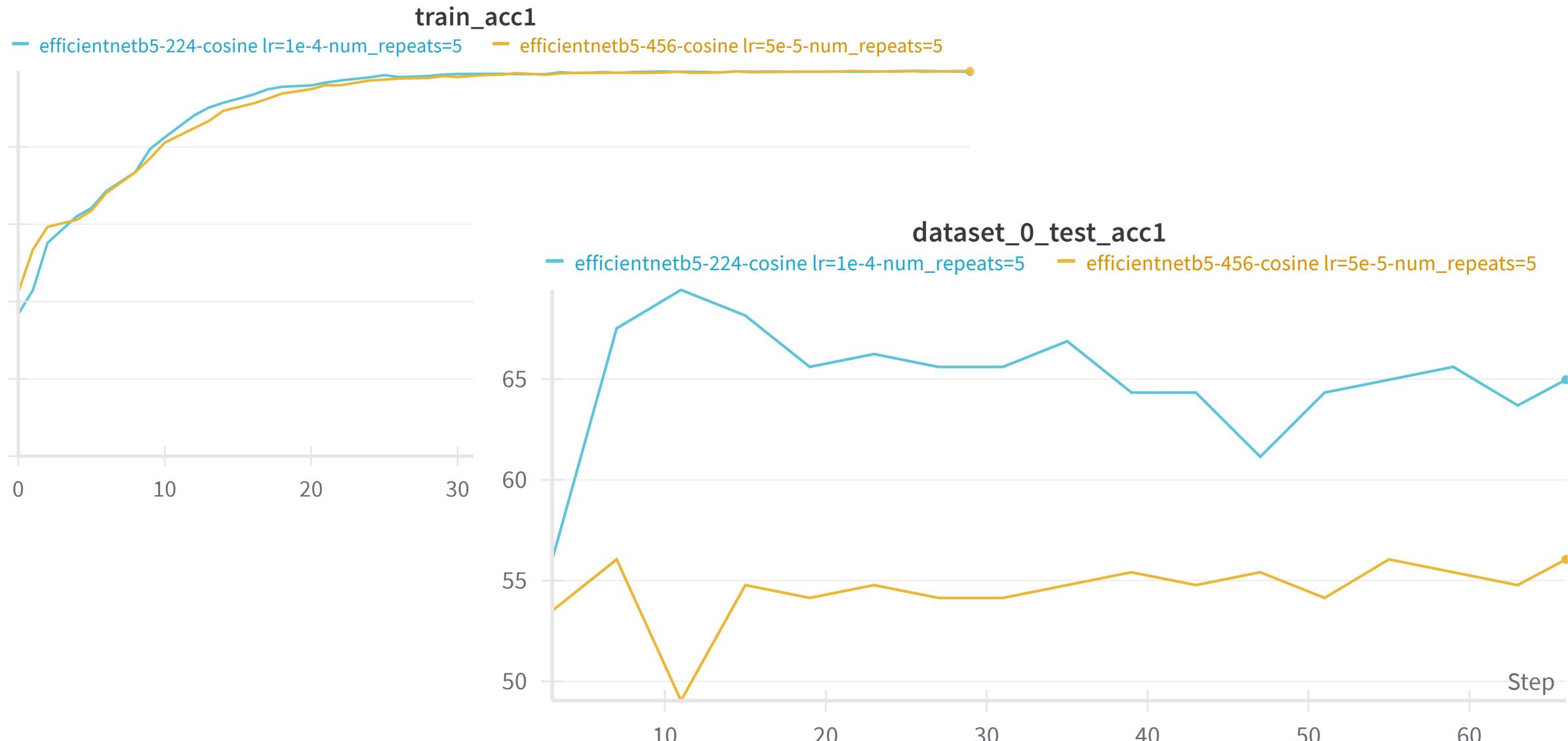
Intuitively, as image resolution increases, accuracy ought to increase as well, since more features are included in the training data and the representation ability of the model also increases.

## Explanation 2: Actually, no need for so many features!



- The medical image data is **Grayscale** and is thus **naturally simpler** than RGB images on which the DL models are pretrained.
- This medical image CLS task **only contains three classes**, so it does not need a large number of latent features for effective classification, unlike typical image CLS tasks where the model needs adequate features so that it can classify an image from typically hundreds of different classes.

# And the Consequence of too many features ... Overfitting!



## Ablation 3: Data Augmentation & Accuracy

|           | #Repeated Augmentation | Top-1 Accuracy |
|-----------|------------------------|----------------|
| Swin_Tiny | 0                      | 78.9           |
| Swin_Tiny | 3                      | <b>80.9</b>    |
| Swin_Tiny | 10                     | 80.3           |

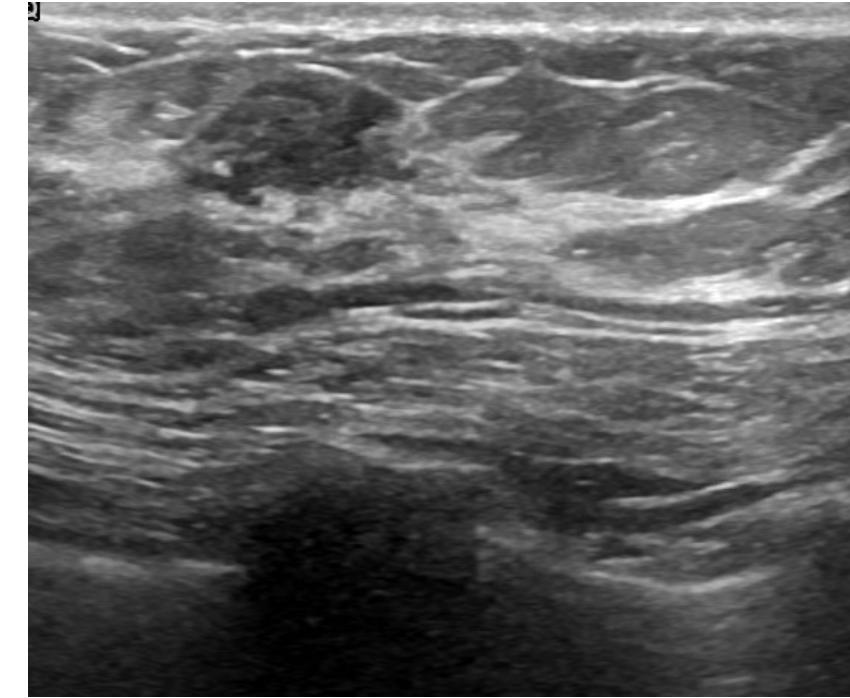
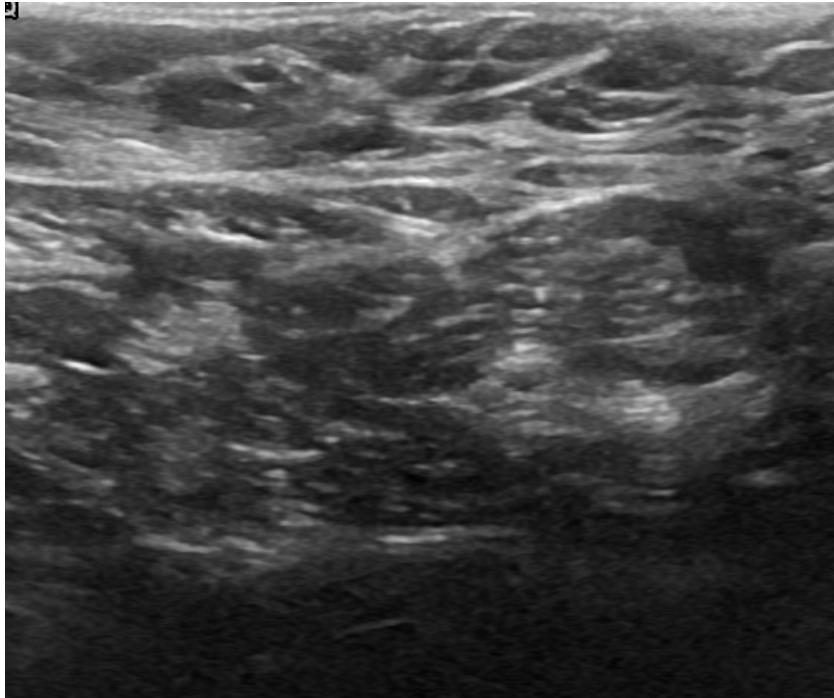
### Analysis:

1. A moderate number of repeated augmentation is helpful to this few-shot image CLS setting, as data augmentation not only **increases the total size of the training set** but also **promotes generalization to a broader range of data**.
2. But why does accuracy saturate and even slightly decrease as the number of repeated augmentation continues to grow?

## Explanation 3: There's a limit in the distribution that data augmentation can fit!

1. First, data augmentation generates new data based on the original training set, so the biases in the original dataset persist in the augmented data. Popular data augmentation methods like Cropping, Flipping, Rotation and ColorJitting **cannot generate completely “new” data that perfectly bridge the gap between the training distribution and testing distribution.**
2. Moreover, the quality of augmented data is usually not guaranteed. **Poor quality data can be generated due to inappropriate perturbation** like overly cropping or weird color jitting, and such data are not guaranteed to be relevant enough to the original training set. Such poor quality or irrelevant data can introduce noise, bias, or inconsistency to the model, leading to inaccurate or misleading predictions.

**Specifically, augmentation on the Grayscale Medical Image Data are especially prone to introduce noise!**



The above two images are of two different classes!

## Ablation 4: Still Model Size, but Between Different Architectures

|                     | #Parameters | Top-1 Accuracy |
|---------------------|-------------|----------------|
| <b>Efficient_b5</b> | <b>30M</b>  | <b>69.4</b>    |
| <b>Swin_Tiny</b>    | <b>29M</b>  | <b>80.9</b>    |

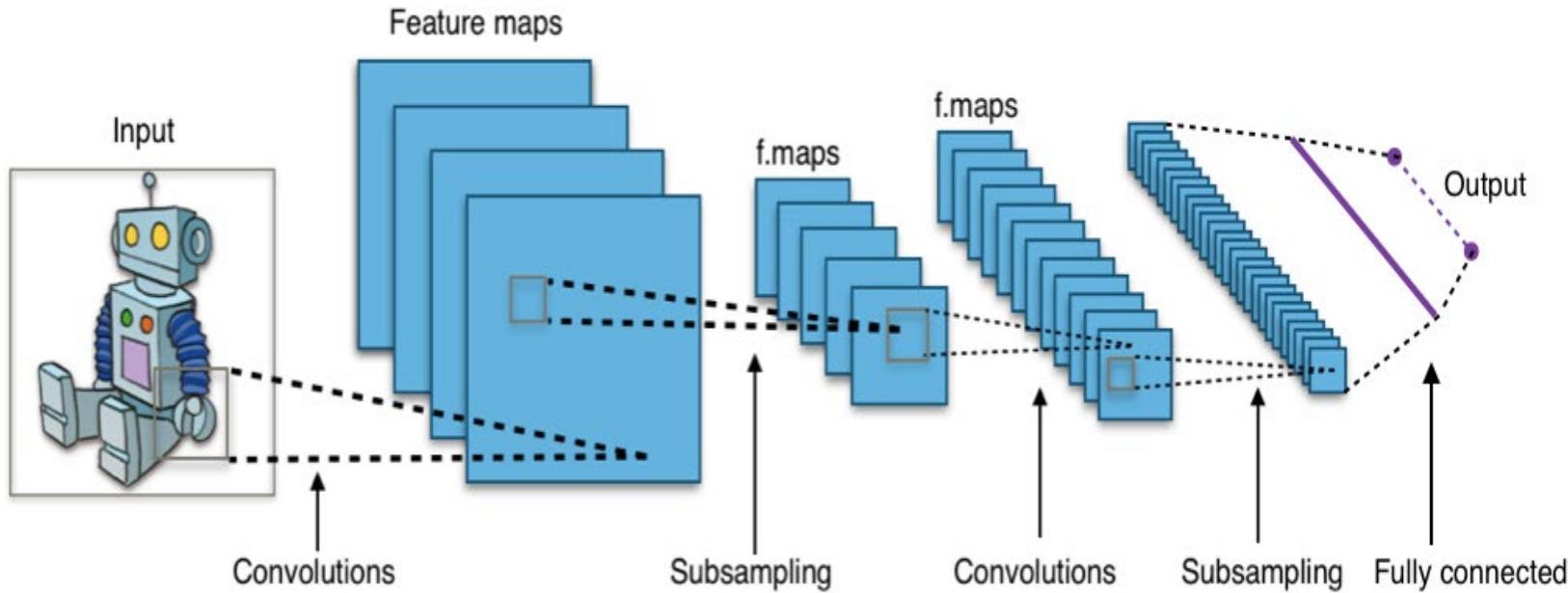
|                     | #Parameters | Top-1 Accuracy |
|---------------------|-------------|----------------|
| <b>Efficient_b7</b> | <b>66M</b>  | <b>67.5</b>    |
| <b>Swin_Base</b>    | <b>88M</b>  | <b>83.4</b>    |

### Counterfactual:

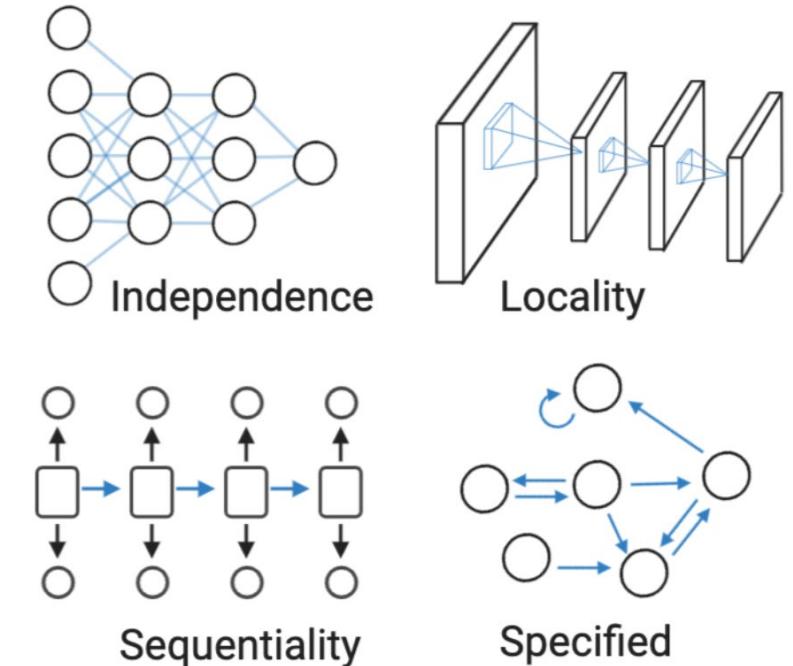
Compare EfficientNets with ViTs of similar parameter amount. The performance of EfficientNets is consistently and significantly poorer than ViTs of similar parameter amount. What's the mystery behind?

# Explanation 4: Inductive Bias - The Inherent and Fundamental Architectural Difference between CNN and ViT

- To summarize, CNNs assume **prior spatial knowledge** in their model design, while ViTs do not have such prior knowledge, and learn every kind of spatial knowledge solely from training and from scratch!
- Such prior knowledge sometimes boosts training under data-scarce scenarios, but more often **constitutes a kind of bias** when the model is already pretrained on a large enough amount of data, because such prior knowledge **obstructs the model from learning long-distance relations between local patterns in an image**.



Relational Inductive Biases



## An Excerpt from UvA DL Notebook

- Convolutional Neural Networks have been designed with the assumption that images are translation invariant. Hence, we apply convolutions with shared filters across the image. Furthermore, *a CNN architecture integrates the concept of distance in an image: two pixels that are close to each other are more related than two distant pixels.* Local patterns are combined into larger patterns until we perform our classification prediction. All those aspects are inductive biases of a CNN.
- In contrast, *a Vision Transformer does not know which two pixels are close to each other, and which are far apart. It has to learn this information solely from the sparse learning signal of the classification task.* This is a huge disadvantage when we have a small dataset since such information is crucial for generalizing to an unseen test dataset. With large enough datasets and/or good pre-training, a Transformer can learn this information without the need for inductive biases, and instead is **more flexible** than a CNN. *Especially long-distance relations between local patterns can be difficult to process in CNNs, while in Transformers, all patches have the distance of one.*

# Summary for Deep Learning Methods

1. Deep Learning Methods, especially ViTs, are significantly more capable in classifying few-shot medical images compared with traditional machine learning methods.
2. Among different deep learning architectures, ViTs are consistently and significantly more capable than CNNs, probably due to their less inductive bias and stronger long-range modelling ability. They show significantly lower degree of overfitting than CNNs.
3. The statistical learning methods are more interpretable than deep learning ones, especially on such tasks as there are only 3 classes so detailed statistical analysis can be conducted. In contrast, though extensive ablation study is done for the deep learning method, much of the observation can only be intuitively interpreted based on empirical experiences of past practitioners.

## References I

-  Al-Dhabyani, Walid et al. (2020). "Dataset of breast ultrasound images". In: *Data in brief* 28, p. 104863.
-  Chi, Mingmin, Rui Feng, and Lorenzo Bruzzone (2008). "Classification of hyperspectral remote-sensing data with primal SVM for small-sized training dataset problem". In: *Advances in space research* 41.11, pp. 1793–1799.
-  Criminisi, A. et al. (2010). "Regression Forests for Efficient Anatomy Detection and Localization in CT Studies". In: pp. 106–117. DOI: [10.1007/978-3-642-18421-5\\_11](https://doi.org/10.1007/978-3-642-18421-5_11).
-  Dinesh, Paidipati and P Kalyanasundaram (2022). "Medical image prediction for diagnosis of breast cancer disease comparing the machine learning algorithms: SVM, KNN, logistic regression, random forest, and decision tree to measure accuracy". In: *ECS Transactions* 107.1, p. 12681.

## References II

-  Kaganov, Helen, A. Ades, and David S. Fraser (2018). "PREOPERATIVE MAGNETIC RESONANCE IMAGING DIAGNOSTIC FEATURES OF UTERINE LEIOMYOSARCOMAS: A SYSTEMATIC REVIEW". In: *International Journal of Technology Assessment in Health Care* 34, pp. 172 –179. DOI: 10.1017/S0266462318000168.
-  Miller, Kimberly D et al. (2022). "Cancer treatment and survivorship statistics, 2022". In: *CA: a cancer journal for clinicians* 72.5, pp. 409–436.
-  Ramesh Kumar, P and A Vijaya (2022). "Naïve Bayes machine learning model for image classification to assess the level of deformation of thin components". In: *Materials Today: Proceedings* 68. 4th International Conference on Advances in Mechanical Engineering, pp. 2265–2274. ISSN: 2214-7853. DOI: <https://doi.org/10.1016/j.matpr.2022.08.489>.

## References III

-  Varoquaux, Gaël and Veronika Cheplygina (2022). "Machine learning for medical imaging: methodological failures and recommendations for the future". In: *NPJ digital medicine* 5.1, p. 48.
-  Webb, Geoffrey I, Eamonn Keogh, and Risto Miikkulainen (2010). "Naïve Bayes.". In: *Encyclopedia of machine learning* 15.1, pp. 713–714.
-  Zaw, Hein Tun, Noppadol Maneerat, and Khin Yadana Win (2019). "Brain tumor detection based on Naïve Bayes Classification". In: *2019 5th International Conference on Engineering, Applied Sciences and Technology (ICEAST)*, pp. 1–4. DOI: 10.1109/ICEAST.2019.8802562.