# What is a foundation model?

# Does domain specific pre-training help?

Fine-tuning requires ~1000 times less training data!!

**Task 1:** Limit the amount of training data

Build competitive classifiers <u>without</u> <u>any training</u>!

**Task 2:** Limit the number of trained parameters

# Automatically annotate datasets by clustering...?



UMAP of epithelium tissue

- benign
- Gleason 3
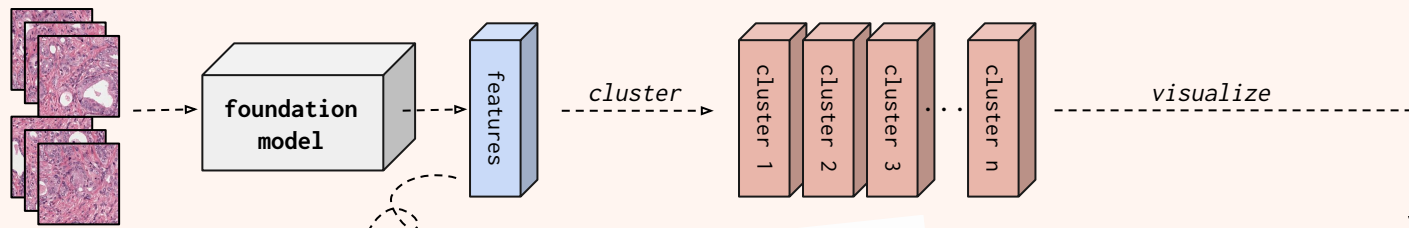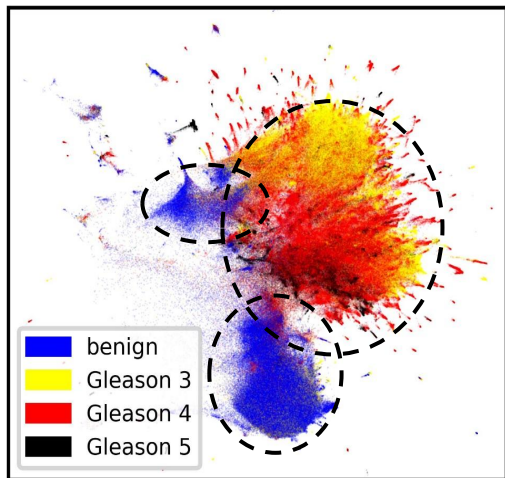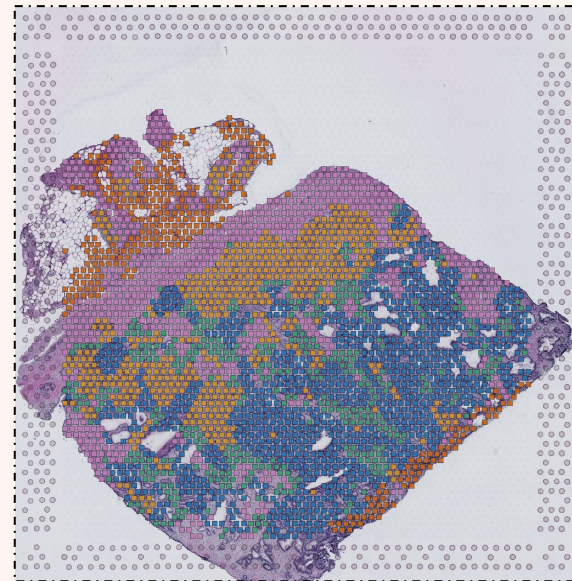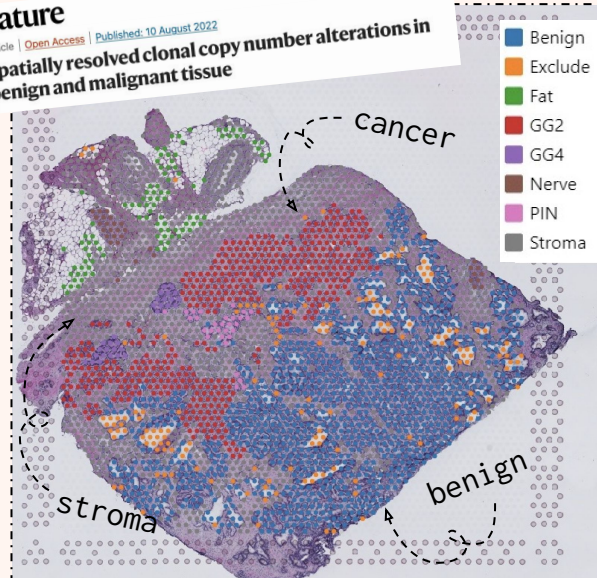- Gleason 4
- Gleason 5
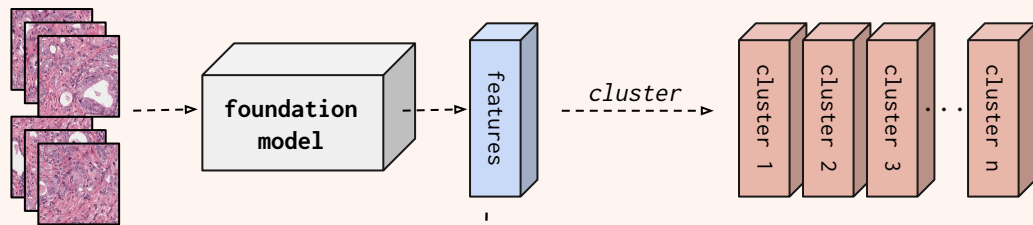
nature

Article | Open Access | Published: 10 August 2022

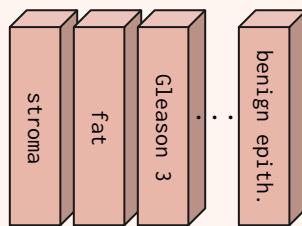Spatially resolved clonal copy number alterations in benign and malignant tissue

cancer

stroma

benign

- Benign
- Exclude
- Fat
- GG2
- GG4
- Nerve
- PIN
- Stroma

# Building multi-modal models...?



Case 1: Tile-level data

Case 2: Patient-level data

Calculate cluster percentages at a **patient-level.**

*patient-level* *cluster percentages*

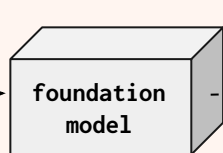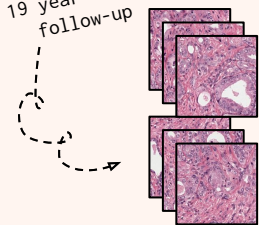| idkey | stroma | fat   | Gleason3 | ... |
|-------|--------|-------|----------|-----|
| 0001  | 0.031  | 0.024 | 0.101    | ... |
| 0002  | 0.183  | 0.210 | 0.017    | ... |
| 0003  | 0.049  | 0.001 | 0.078    | ... |
| ...   | ...    | ...   | ...      | ... |
| 0999  | 0.227  | 0.012 | 0.279    | ... |

Each cluster contains similar histological patterns!

Distribution of histological patterns for a given patient!!
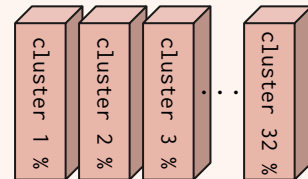
# Prostate cancer specific death

# Summary

Foundation model fine-tuning requires <u>no training</u>
or <u>1000x less training data</u>.

Images with similar histological patterns
produce similar embedded features.

Pre-annotate
whole datasets!!

Combine histomics
with other data
modalities!!!

```
# Cut slide images into small tiles.
HistoPrep -i 'slides/*.tiff' -o tiles/ --width 512
# Extract & cluster features for all tiles.
HistoEncoder extract -i tiles/ -m prostate_medium
HistoEncoder cluster -i tiles/ -n 8 16 32 64 128
```

https://github.com/jopo666/HistoEncoder