# Cancer Cell Cluster Detection by Unified Tumor Tissue Slide Image and Gene Expression Quantification

**Wenkang Zhan, Zhengkang Fan**
**Dept. of Computer and Information Science**

# Background

Diagnosis of cancer region in a slide image is difficult and time-consuming, as cancer may transfer to other place.

Machine learning have many successful applications on image data. We want to use techniques to detect cancer cluster region.

In real scenarios, slide image data set is small. Gene expression quantification is available and it is a powerful material to classify tumors. Thus, we unify gene expression quantification and slide image to improve the performance of cancer cell cluster detection.
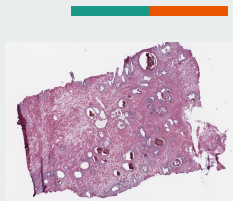
# Data Collection

Data source: The Cancer Genome Atlas (TCGA)
- Public-open database
- Provide over 20,000 cancers and matches 33 cancer types.
- https://portal.gdc.cancer.gov/

Data:
- Slide images
- Gene expression quantification -> feature
- Clinical data: diagnosis record -> label
- Meta data: Identity information ->tag for combining data
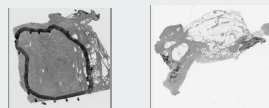- supplemental clinical data: unprocessed data containing some useful information

Data Collection and Pre-processing

# Feature Reduction - gene expression quantification



Correlation of feature v.s. label

Set threshold to select features

Gene expression quantification
(over 60,000 features)

Manually select meaningful features

Performance comparison with different subset of features by a baseline models

Gene expression quantification
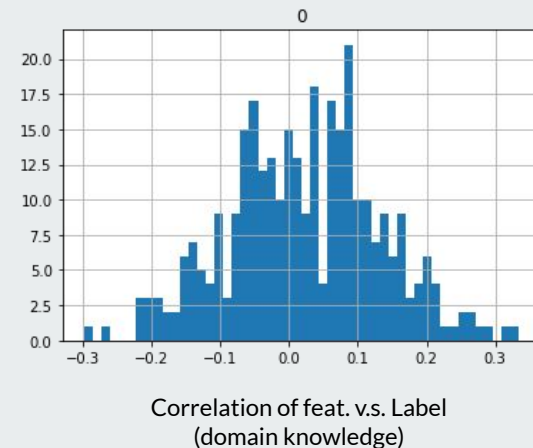(hundreds of features)

# Feature Reduction - gene expression quantification

**Number of remaining features with different thresholds**

| threshold | 0.16 | 0.18 | 0.2 | 0.22 | 0.24 | 0.26 | 0.28 | 0.30 |
|---|---|---|---|---|---|---|---|---|
| #features | 3566 | 2333 | 1534 | 969 | 623 | 397 | 238 | 132 |

**Accuracy of baselines - only on gene expression quantification**

| | Selected features by correlations | | | | | | | | Domain knowledge |
|---|---|---|---|---|---|---|---|---|---|
| thres. | 0.16 | 0.18 | 0.20 | 0.22 | 0.24 | 0.26 | 0.28 | 0.30 | - |
| #feat. | 3566 | 2333 | 1534 | 969 | 623 | 397 | 238 | 132 | 315 |
| LR | **0.66** | 0.66 | 0.6 | 0.63 | 0.63 | 0.6 | 0.56 | 0.56 | **0.66** |
| CNN | 0.44 | 0.43 | 0.43 | 0.48 | 0.48 | **0.56** | **0.56** | **0.56** | **0.56** |
| LSTM | 0.46 | 0.48 | 0.48 | 0.48 | 0.48 | 0.53 | **0.56** | **0.56** | **0.56** |



Correlation of feat. v.s. Label
(domain knowledge)

# Model



Feature Selected Gene expression data

Regression: Y=XW+B

Cropped image data

Input layer  Conv2d  Max pooling  Conv2d  Dense layer

3 conv. layers in total

Dense layer

Prediction
（Pathological stage）

Marked image data

# Results

# Results



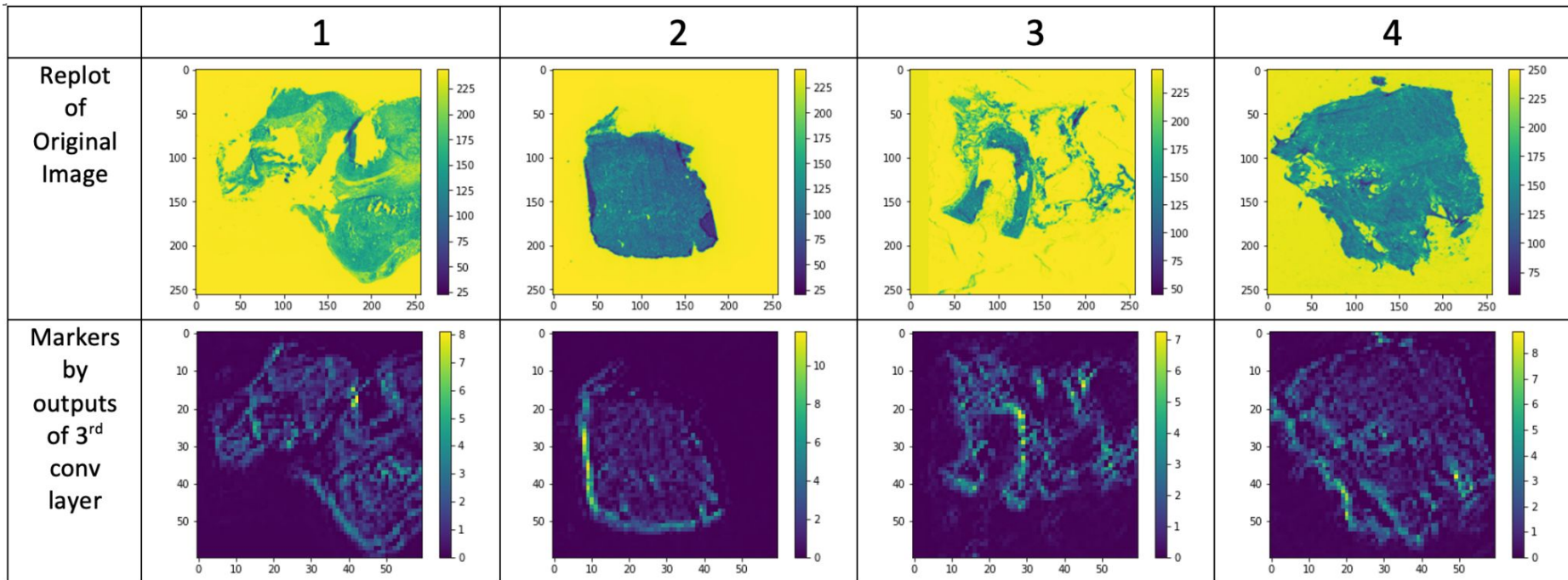|  | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Replot of Original Image | | | | |
| Markers by outputs of 3<sup>rd</sup> conv layer | | | | |

# Future

- Add more data into dataset(115 samples)

- Overfitting

- Improve the model structure

# Thanks for your listening!

# Q&A

**Wenkang Zhan, Zhengkang Fan**
**Dept. of Computer and Information Science**