

## Motivation

### Multimodal Pathology [1]

- Patient-level outcome prediction from fusion of complementary biological information – e.g.) Tissue **Morphology** + **Transcriptomics** expression

### Morphology

- **Digitized tissue sections** (whole-slide images, WSIs), of up to **100,000 x 100,000** pixels (at 0.5μm/pixel)
- Typically tokenized into > 10,000 patch tokens (256 x 256 pixels)

### Transcriptomics (Genes)

- Whole-Transcriptome RNA-sequencing provides expressions for > 20,000 genes (tokens)

### Limitations of current multimodal approaches

- Multimodal fusion of large sets of token embeddings often leads to
  - Computationally-infeasible training
  - Unstable training dynamics for survival prediction

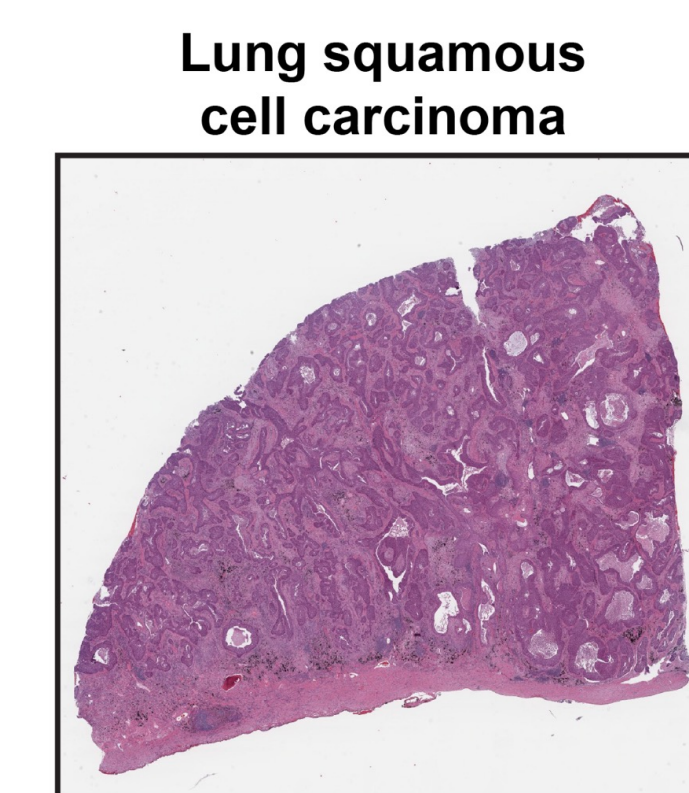
Can we create a **token-efficient multimodal framework**?

## Multimodal patient representation

### Morphology

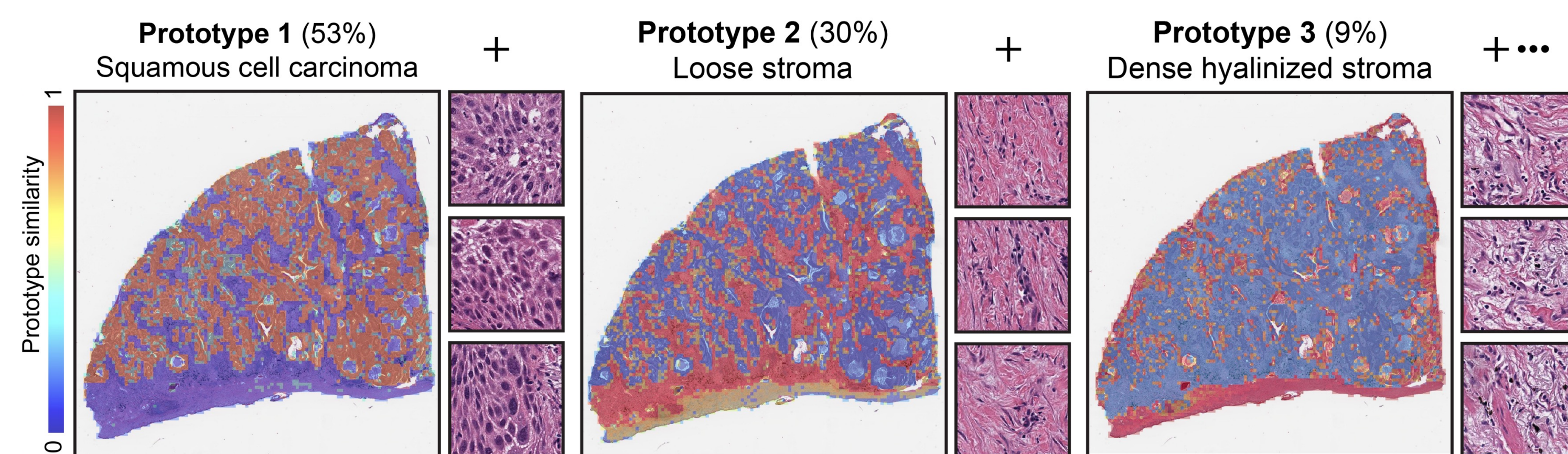
#### Redundant morphological information in WSI

- Handful of morphological patterns repeated throughout the tissue (e.g., cancer cells, stroma, adipose tissue)



#### Prototype-based summarization of WSI

- **WSI**  $\cong$  Distribution of **morphological concepts**
- Summarization of WSI based on two important conditions
  - Feature representation of each concept
  - Cardinality (proportion) of each concept in WSI
- Huge compression
  - Prototypes ( $C = 8 \sim 32$ )  $\ll$  patches per WSI ( $N \sim 10^4$ )
- Optimal transport or Gaussian mixture models [2]



### Transcriptomics

Only sparse set of genes are relevant for cancer

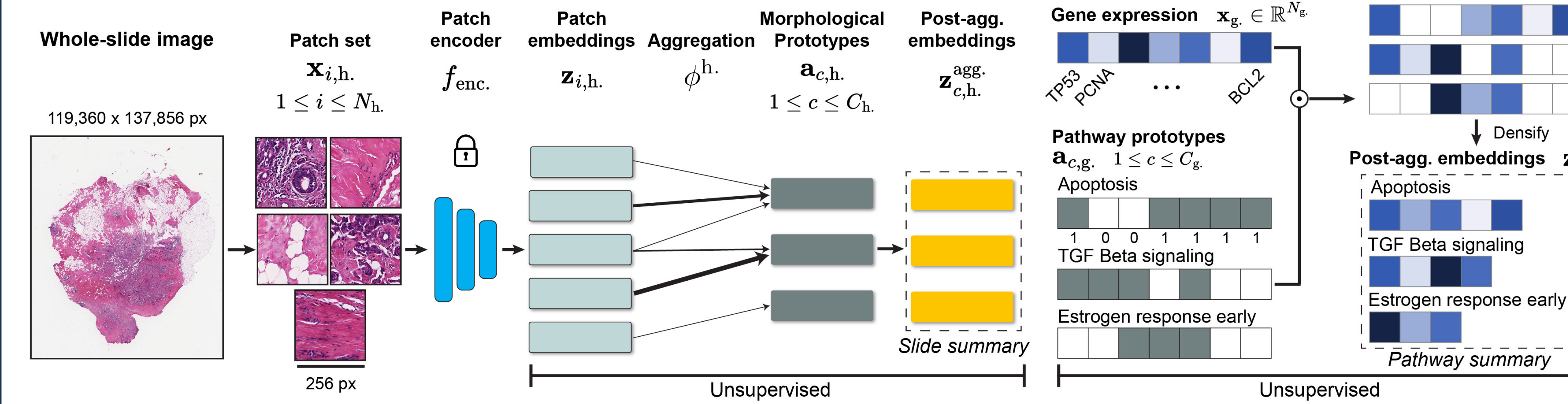
#### Prototype-based summarization (pathway) of transcriptomics

- **Pathways** provide natural functional grouping of genes
  - **Example Pathway:** HALLMARK\_APOPTOSIS
  - **Associated genes** (161 genes): [ADD1, AIFM3, ..., WEE1, XIAP]
- **Pathways** can be effectively treated as **prototypes**
  - **Semantic group:** Apoptosis = Programmed cell death
  - **Compression:** > 20,000 genes  $\Rightarrow$  A functional group of 200 genes

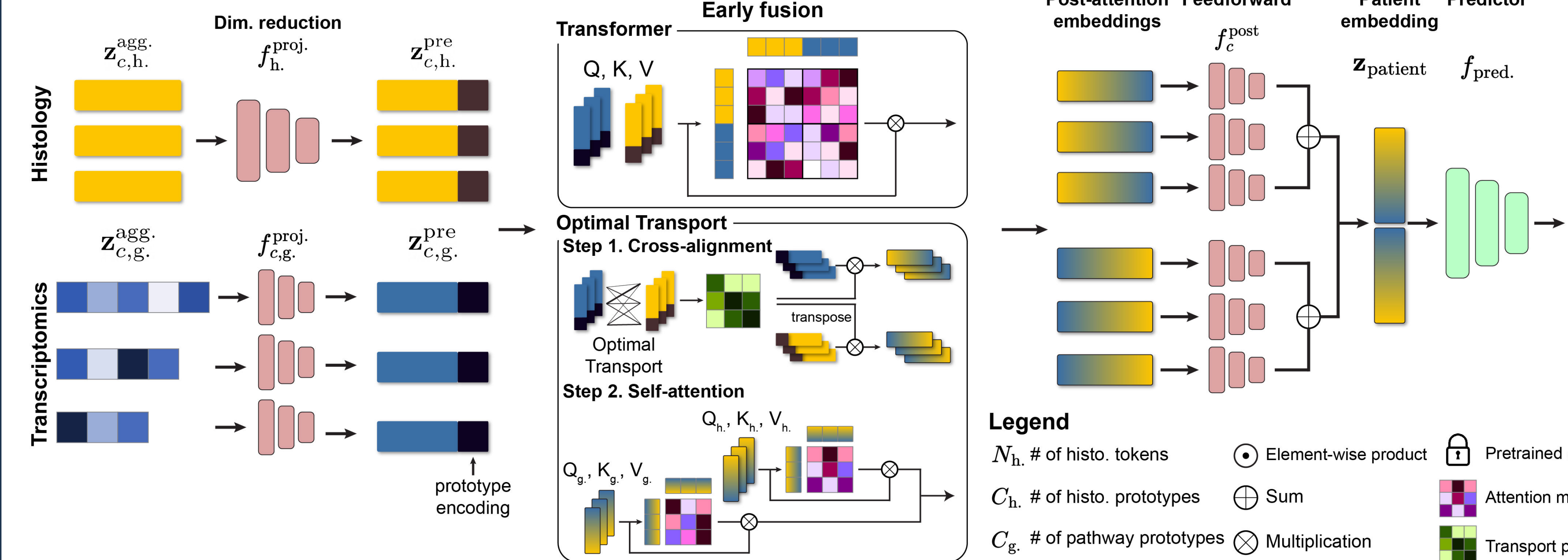
$\Rightarrow$  **MMP: MultiModal Prototyping** for survival prediction

## MMP for multimodal cancer survival prediction

### A. Morphological prototypes



### C. Multimodal fusion



### Morphology: Patch embedding distribution (Gaussian Mixture Model)

- $p(\mathbf{z}_n; \theta) = \sum_{c=1}^C \pi_c \cdot N(\mathbf{z}_n; \mu_c, \Sigma_c) \Rightarrow$  Each component: a prototype and its distribution
- **Morphology prototype set:**  $\mathbf{z}_{\text{histo}}^{\text{pre}} = \{\mathbf{z}_{c, \text{histo}}^{\text{pre}}\}_{c=1}^C \Rightarrow \mathbf{z}_{c, \text{histo}}^{\text{pre}} = f([\hat{\pi}_c, \hat{\mu}_c, \hat{\Sigma}_c]) \in \mathbb{R}^d$

### Transcriptomics: 50 Hallmark functional gene sets (Pathways)

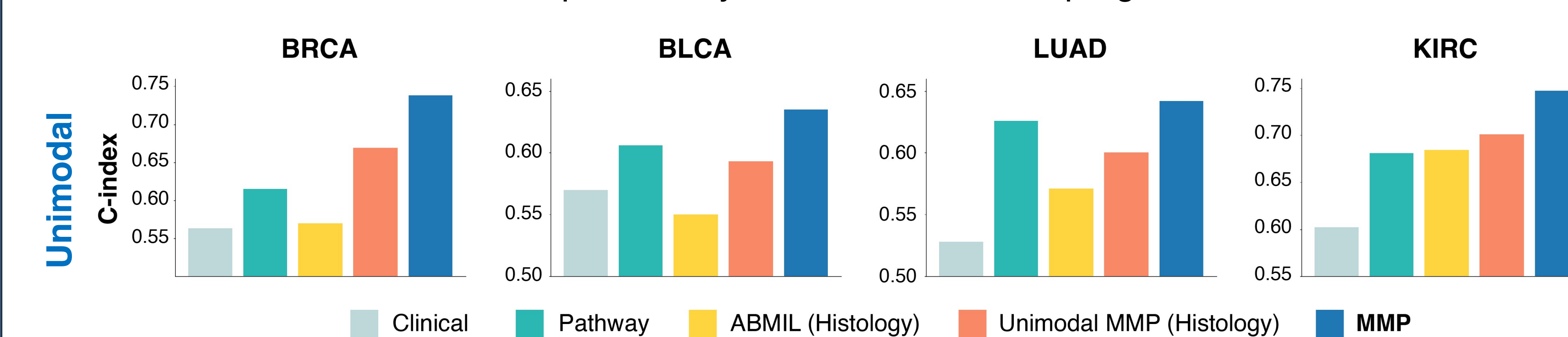
- **Pathway prototype set:**  $\mathbf{z}_{\text{gene}}^{\text{pre}} = \{\mathbf{z}_{c, \text{gene}}^{\text{pre}}\}_{c=1}^{50} \Rightarrow \mathbf{z}_{c, \text{gene}}^{\text{pre}} \in \mathbb{R}^d$

### Multimodal early fusion (Morphology prototypes + Pathway prototypes)

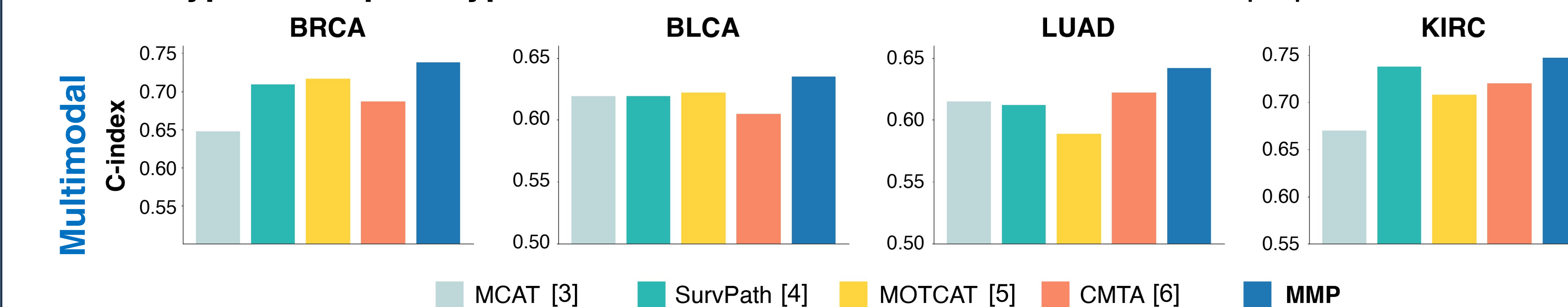
- Transformer-based fusion or Optimal Transport-based fusion

## MMP for slide-level survival prediction

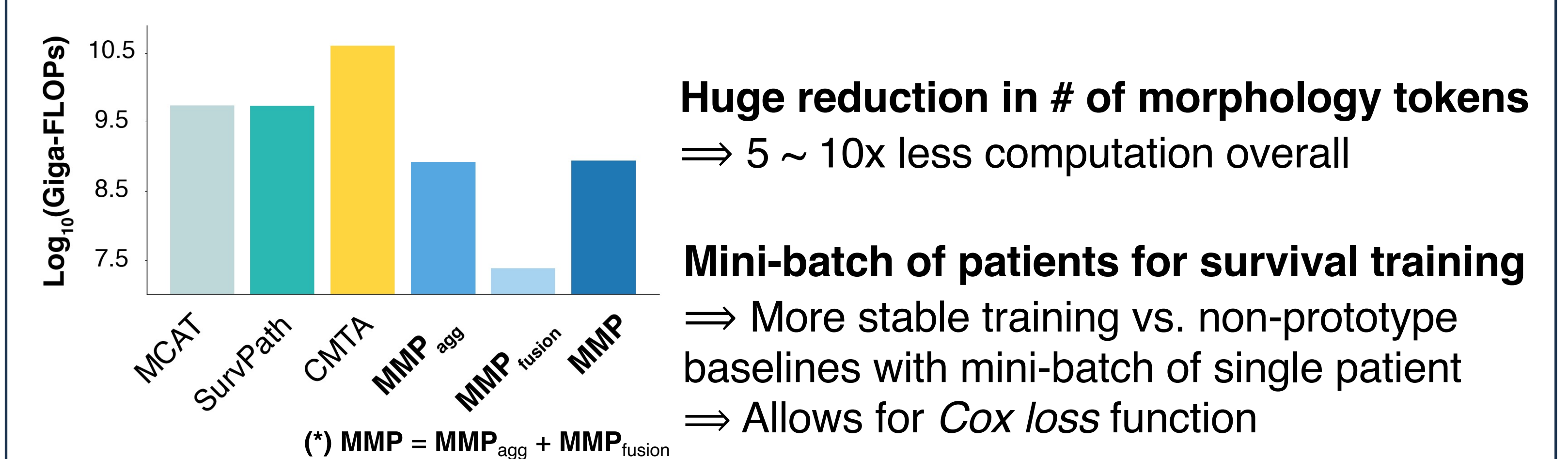
- **Multimodal > Unimodal:** Complementary information benefits prognosis



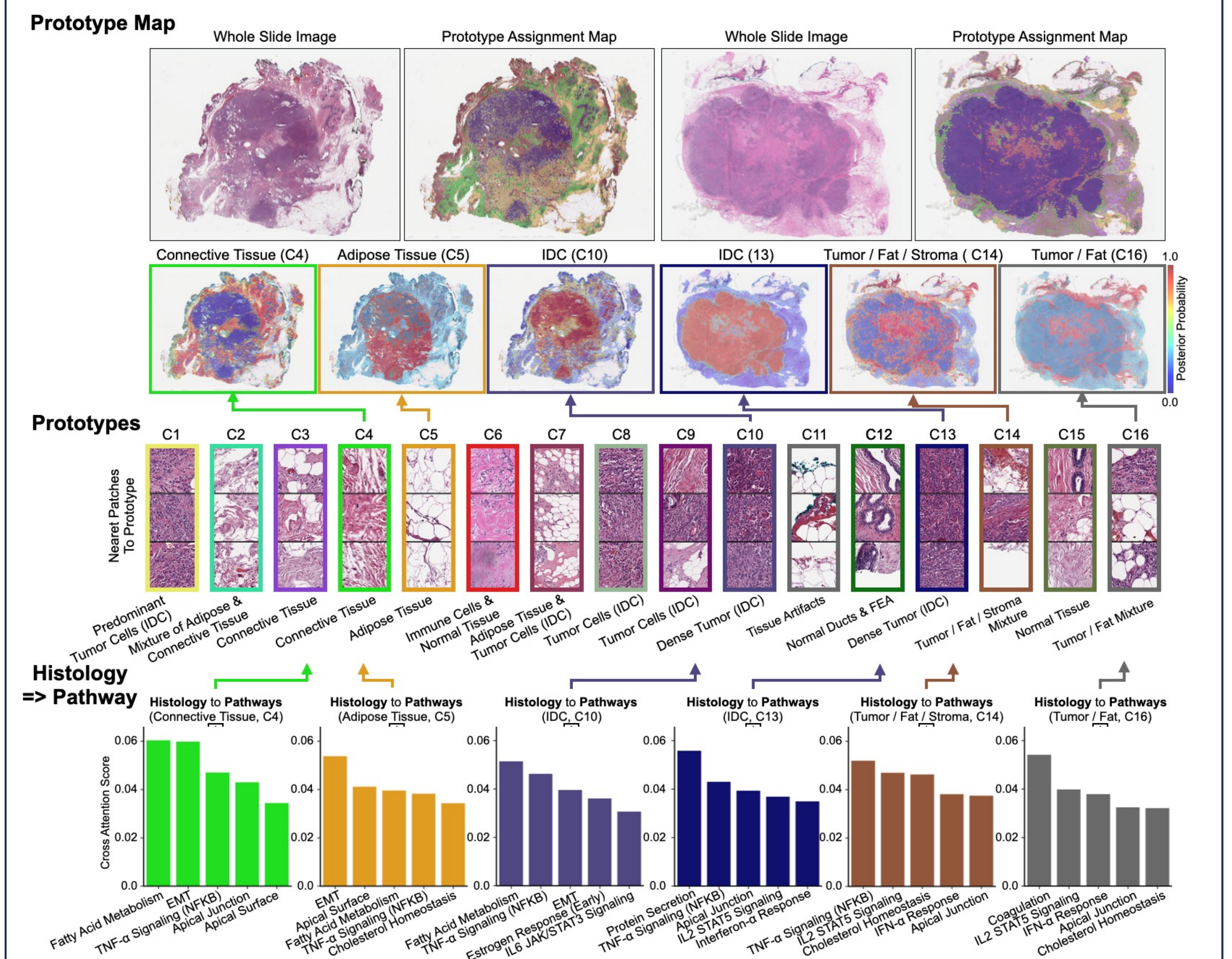
- **Prototype > non-prototype multimodal baselines:** Smaller token set helps performance



## MMP computational complexity



## MMP for Interpretability



### Prototype-oriented interpretability

- Visualization of the most similar prototype on WSI (Cluster map)
- Each prototype represents different morphological concepts

### Cross-modal interpretability

- Bidirectional (Histology  $\rightarrow$  Pathway, Pathway  $\rightarrow$  Histology)
- Intuitive interpretability based on prototypes in both domains

## References

- [1] Song AH et al., Artificial intelligence for digital and computational pathology. *Nature Reviews Bioengineering*, 2023
- [2] Song AH et al., Morphological prototyping for unsupervised slide representation learning in computational pathology. *CVPR*, 2024
- [3] Chen RJ et al., Multimodal co-attention transformer for survival prediction in gigapixel whole slide images, *ICCV*, 2021
- [4] Jaume G et al., Modeling dense multimodal interactions between biological pathways and histology for survival prediction, *CVPR*, 2024
- [5] Xu Y et al., Multimodal optimal transport-based co-attention transformer with global consistency for survival prediction, *ICCV*, 2023
- [6] Zhou F et al., Cross-modal translation and alignment for survival analysis, *ICCV*, 2023

PAPER



CODE

