

Prior Knowledge Injection into Deep Learning Models Predicting Gene Expression from Whole Slide Images

Max Hallemeesch, Marija Pizurica, Paloma Rabaey, Olivier Gevaert, Thomas Demeester, Kathleen Marchal

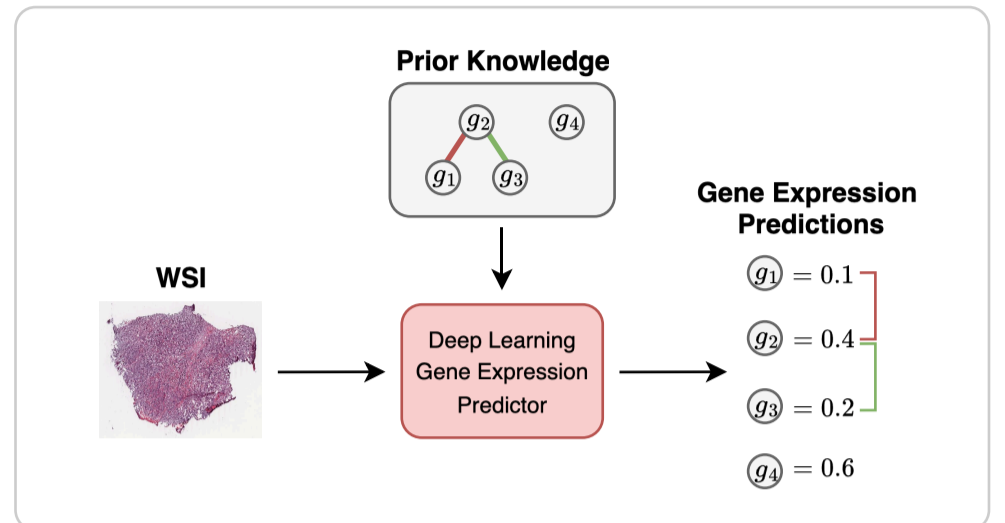
Problem Statement

Recent advances in Deep Learning allow to predict molecular information from morphological features within Whole Slide Images (WSIs). While promising, **current methods lack the robustness** to fully replace direct sequencing.

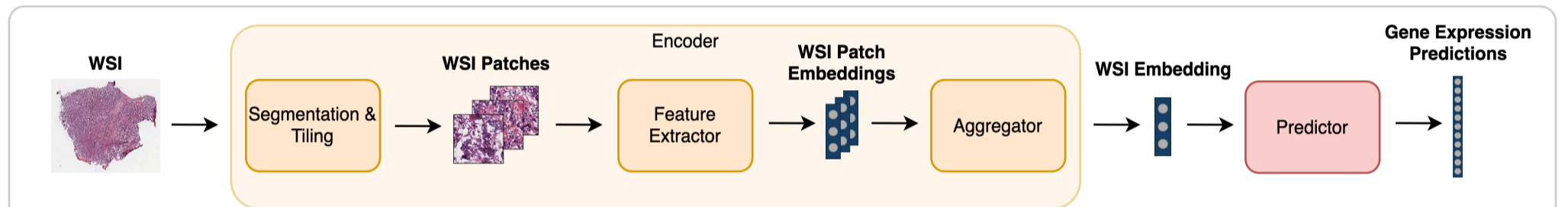
Goal

Here we aim to improve existing methods by introducing a **model-agnostic framework that allows to inject prior knowledge on gene-gene interactions** into Deep Learning architectures.

High Level Approach

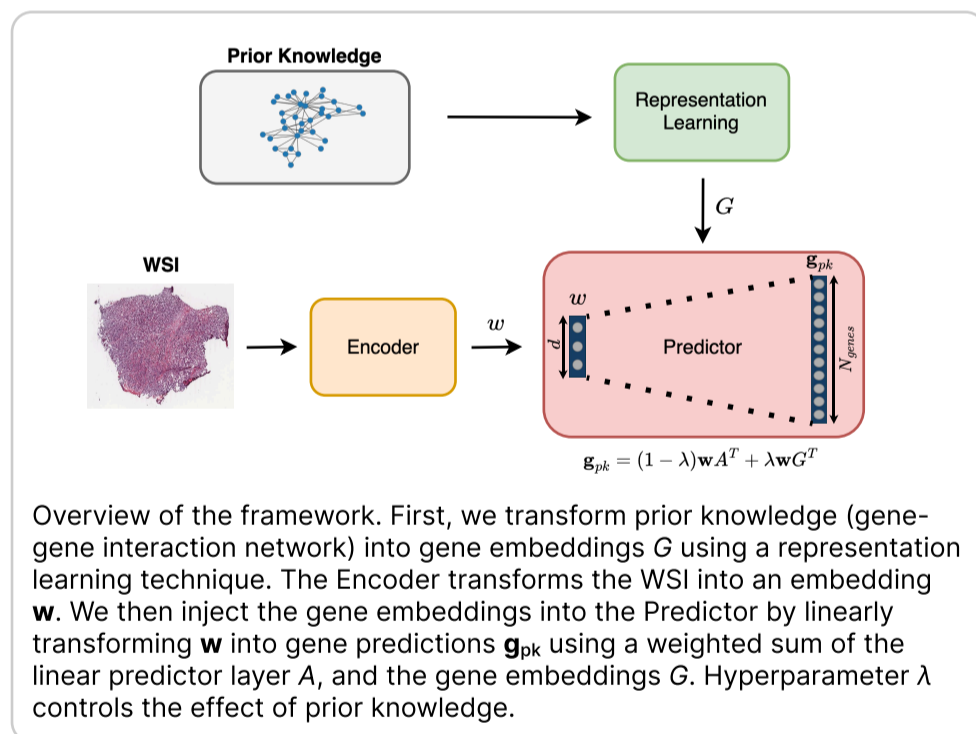


General Workflow Gene Expression Prediction

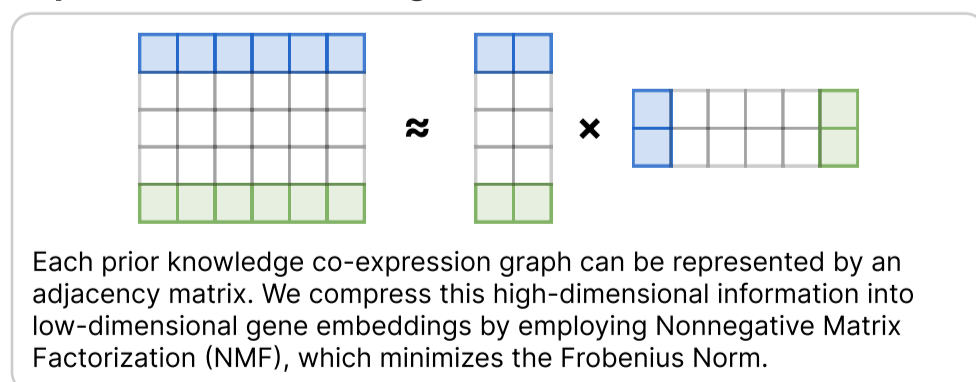


General workflow of predicting gene expression from Whole Slide Images. First, the WSI is processed by an encoder, which extracts patches, their corresponding features, and aggregates them into a single WSI embedding. Then, a predictor transforms the embedding into expression predictions for 25,761 genes. We evaluate our model-agnostic framework by considering two feature extractors (CTrans (ctr) and UNI) and three aggregators (MLP, Transformer (tf) and SummaryMixing (smx)).

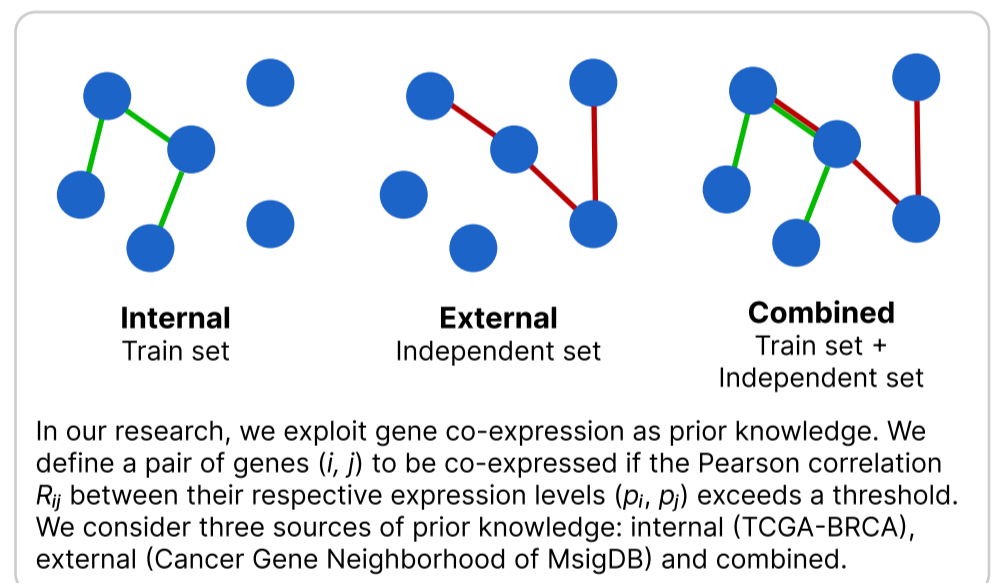
Framework Architecture



Representation Learning



Prior Knowledge Sources



Evaluation

TCGA	No PK	External	Internal	Combined
ctr_mlp	21,233	22,225 (0.9) $\uparrow 992$	22,278 (0.9) $\uparrow 1,045$	22,160 (0.9) $\uparrow 927$
ctr_tf	19,155	21,647 (0.9) $\uparrow 2,429$	20,618 (0.2) $\uparrow 1,463$	20,548 (0.8) $\uparrow 1,393$
ctr_smx	20,945	22,564 (0.9) $\uparrow 1,619$	21,451 (0.5) $\uparrow 506$	21,944 (0.5) $\uparrow 999$
uni_mlp	21,721	22,666 (0.8) $\uparrow 945$	22,802 (0.8) $\uparrow 1,081$	23,214 (0.9) $\uparrow 1,493$
uni_tf	22,124	22,461 (0.2) $\uparrow 337$	22,645 (0.1) $\uparrow 521$	22,597 (0.1) $\uparrow 473$
uni_smx	22,997	23,578 (0.5) $\uparrow 581$	23,732 (0.1) $\uparrow 735$	23,162 (0.9) $\uparrow 165$

CPTAC	No PK	External	Internal	Combined
ctr_mlp	16,936	16,313 $\downarrow 623$	18,116 $\uparrow 1,180$	17,363 $\uparrow 427$
ctr_tf	15,677	15,146 $\downarrow 531$	14,983 $\downarrow 694$	14,386 $\downarrow 1,291$
ctr_smx	15,714	16,682 $\uparrow 968$	15,731 $\uparrow 17$	15,753 $\uparrow 39$
uni_mlp	15,560	16,106 $\uparrow 546$	16,045 $\uparrow 485$	15,784 $\uparrow 224$
uni_tf	14,705	15,648 $\uparrow 763$	15,469 $\uparrow 764$	15,400 $\uparrow 695$
uni_smx	16,952	17,280 $\uparrow 328$	17,091 $\uparrow 139$	16,981 $\uparrow 29$

We evaluated the number of significantly predicted genes across 18 experiments, including three sources of prior knowledge and six deep learning architectures, on both TCGA-BRCA (upper table) and CPTAC-BRCA (lower table). Across 14 experiments we observed an increase in the number of significant genes on both TCGA and CPTAC, demonstrating an enhanced generalization performance.

