# Multimodal Lego: Model Merging and Fine-Tuning Across Topologies and Modalities

UNIVERSITY OF CAMBRIDGE — GATES Cambridge

Konstantin Hemker, Nikola Simidjievski, Mateja Jamnik    {kh701, ns779, mj201}@cam.ac.uk

Code    Paper    ICLR

## Problem

**Multimodal** models require many **paired training samples** for competitive performance.

**Fusion & ensembling** methods may sacrifice performance from **signal interference**

**Model merging** requires equivalent architectures to interpolate weights.

How can we build performant multimodal models from pre-trained unimodal encoders?
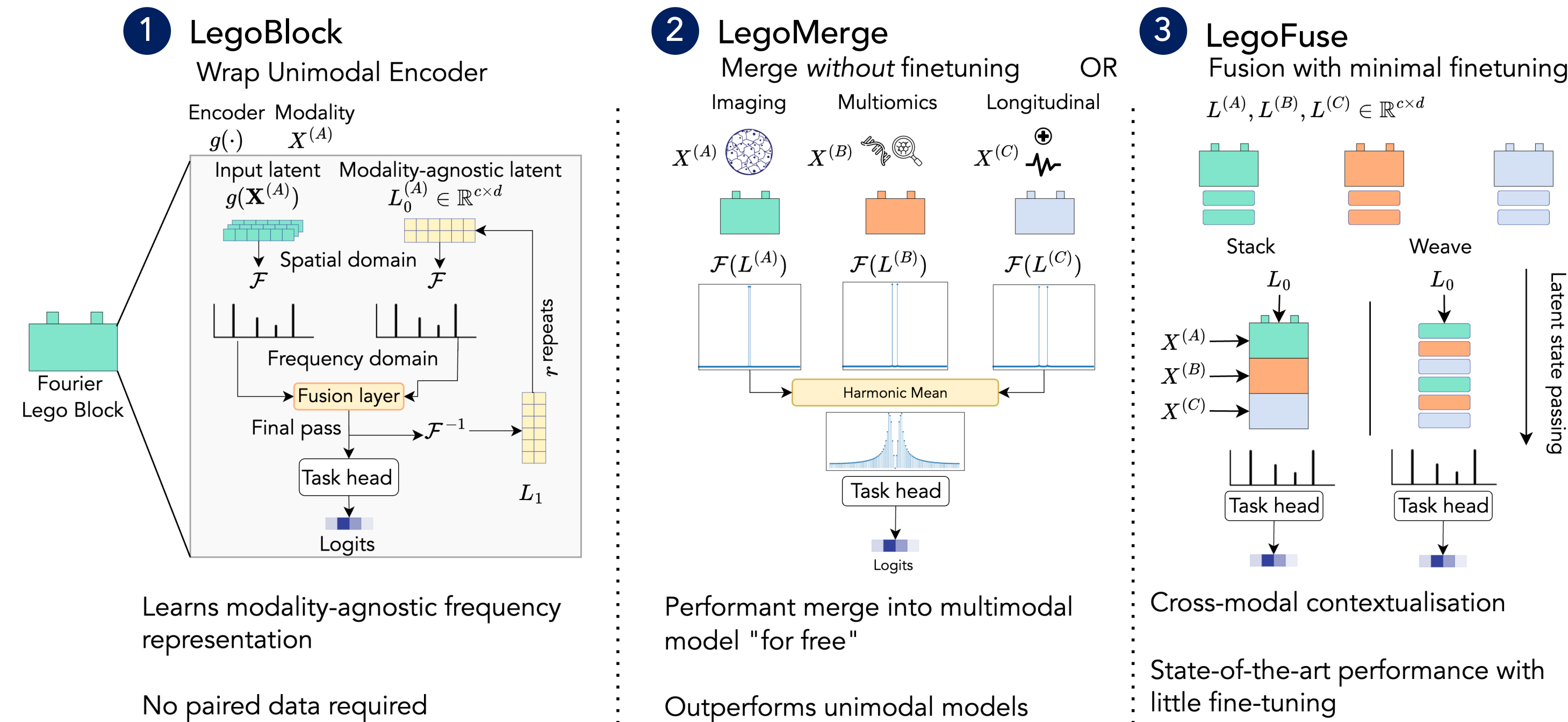
## Solution

**Novel multimodal model merging paradigm** with three components:

1 **LegoBlocks:** fits modality-specific adapter to pre-trained model with any topology

2 **LegoMerge:** effectively merges the blocks with little signal interference

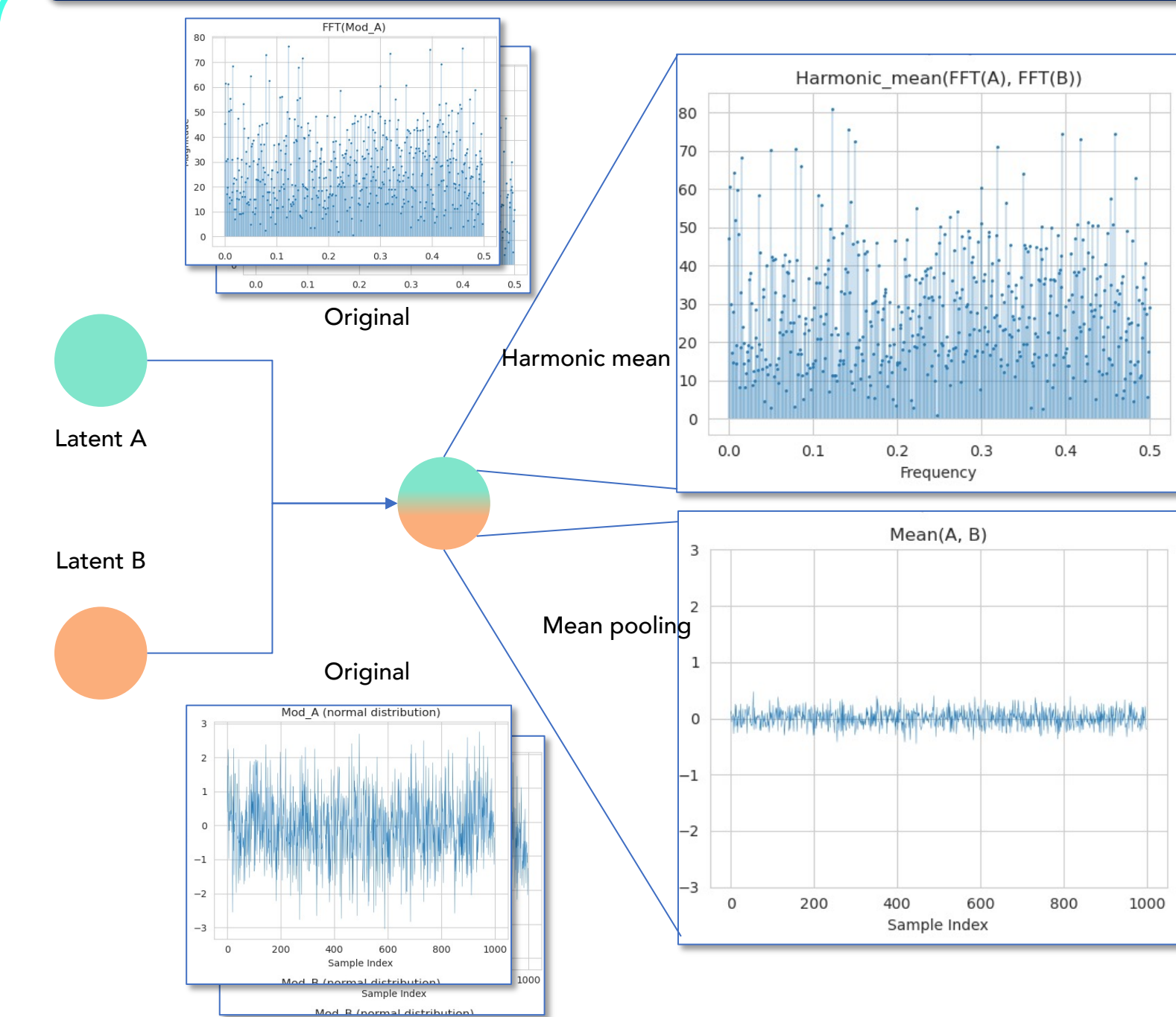3 **LegoFuse:** allows for parameter-efficient fine-tuning

## Contributions

✓ Performant multimodal models without costly e2e training

✓ Agnostic to model architecture, enabling flexible multimodal learning

✓ Scalable to any number of modalities

✓ Robust to high cross-modal imbalance and/or missing modalities

### Architecture Overview



1 **LegoBlock** — Wrap Unimodal Encoder
Learns modality-agnostic frequency representation
No paired data required

2 **LegoMerge** — Merge *without* finetuning OR
Imaging · Multiomics · Longitudinal
Performant merge into multimodal model "for free"
Outperforms unimodal models

3 **LegoFuse** — Fusion with minimal finetuning
$L^{(A)}, L^{(B)}, L^{(C)} \in \mathbb{R}^{c \times d}$
Stack · Weave
Cross-modal contextualisation
State-of-the-art performance with little fine-tuning

### Frequency-domain representations preserve signal



### Results

| | BLCA | BRCA | KIRP | UCEC | ICD9 | MORT | ISIC |
|---|---|---|---|---|---|---|---|
| *Samples* | n=436 | N=1021 | n=284 | n=538 | n=32616 | n=32616 | n=2875 |
| Modalities | img, mut, cnv, rna | img, mut, cnv, rna | img, mut, cnv, rna | img, mut, cnv, rna | tab, ts | tab, ts | tab, img |
| *Metric* | c-Index | c-Index | c-Index | c-Index | AUC | Macro AUC | AUC |
| **Unimodal (Tabular)** | | | | | | | |
| SNN | $0.689_{\pm0.012}$ | $0.544_{\pm0.020}$ | $0.798_{\pm0.035}$ | $0.589_{\pm0.057}$ | $0.731_{\pm0.023}$ | $0.634_{\pm0.020}$ | $0.507_{\pm0.005}$ |
| MultiModN | $0.500_{\pm0.000}$ | $0.500_{\pm0.000}$ | $0.525_{\pm0.140}$ | $0.500_{\pm0.000}$ | $0.500_{\pm0.000}$ | $0.500_{\pm0.000}$ | $0.500_{\pm0.000}$ |
| Perceiver | $0.686_{\pm0.009}$ | $0.557_{\pm0.016}$ | $0.836_{\pm0.053}$ | $0.615_{\pm0.035}$ | $0.629_{\pm0.023}$ | $0.658_{\pm0.009}$ | $0.840_{\pm0.084}$ |
| *LegoBlock* | $0.681_{\pm0.015}$ | $0.591_{\pm0.021}$ | $0.840_{\pm0.135}$ | $0.615_{\pm0.031}$ | $0.645_{\pm0.017}$ | $0.619_{\pm0.028}$ | $0.668_{\pm0.141}$ |
| **Unimodal (Image/Time-series)** | | | | | | | |
| ABMIL | $0.591_{\pm0.057}$ | $0.610_{\pm0.093}$ | $0.741_{\pm0.080}$ | $0.558_{\pm0.040}$ | $0.614_{\pm0.025}$ | $0.691_{\pm0.014}$ | $0.500_{\pm0.000}$ |
| MultiModN | $0.520_{\pm0.022}$ | $0.527_{\pm0.150}$ | $0.570_{\pm0.156}$ | $0.564_{\pm0.097}$ | $0.500_{\pm0.000}$ | $0.544_{\pm0.033}$ | $0.500_{\pm0.000}$ |
| Perceiver | $0.532_{\pm0.027}$ | $0.604_{\pm0.064}$ | $0.716_{\pm0.063}$ | $0.534_{\pm0.106}$ | $0.700_{\pm0.013}$ | $0.715_{\pm0.016}$ | $0.719_{\pm0.050}$ |
| *LegoBlock* | $0.568_{\pm0.029}$ | $0.533_{\pm0.000}$ | $0.630_{\pm0.182}$ | $0.565_{\pm0.069}$ | $0.643_{\pm0.013}$ | $0.711_{\pm0.008}$ | $0.706_{\pm0.147}$ |
| **Multimodal** | | | | | | | |
| *LegoMerge* | $0.701_{\pm0.021}$ | $0.601_{\pm0.025}$ | $0.825_{\pm0.114}$ | $0.625_{\pm0.080}$ | $0.684_{\pm0.015}$ | $0.751_{\pm0.027}$ | $0.721_{\pm0.143}$ |
| *Uplift (Merge vs. best Block)* | 2.9% | 1.7% | -1.8% | 1.6% | 5.7% | 5.3% | 2.1% |
| SNN + ABMIL (CC, Late) | $0.561_{\pm0.000}$ | $0.541_{\pm0.014}$ | $0.841_{\pm0.128}$ | $0.601_{\pm0.018}$ | $0.628_{\pm0.020}$ | $0.617_{\pm0.015}$ | $0.661_{\pm0.196}$ |
| SNN + ABMIL (BL, Late) | $0.622_{\pm0.054}$ | $0.557_{\pm0.089}$ | $0.811_{\pm0.108}$ | $0.666_{\pm0.031}$ | $0.500_{\pm0.000}$ | $0.500_{\pm0.001}$ | $0.501_{\pm0.002}$ |
| Perceiver (CC, Early) | $0.547_{\pm0.060}$ | $0.561_{\pm0.105}$ | $0.692_{\pm0.000}$ | $0.548_{\pm0.000}$ | $0.733_{\pm0.028}$ | $0.723_{\pm0.015}$ | $0.721_{\pm0.198}$ |
| MultiModN (Inter.) | $0.524_{\pm0.018}$ | $0.500_{\pm0.000}$ | $0.602_{\pm0.076}$ | $0.512_{\pm0.006}$ | $0.500_{\pm0.000}$ | $0.500_{\pm0.000}$ | $0.500_{\pm0.000}$ |
| MCAT (Inter.) | $0.702_{\pm0.032}$ | $0.564_{\pm0.000}$ | $0.823_{\pm0.076}$ | $0.633_{\pm0.068}$ | $0.500_{\pm0.000}$ | $0.500_{\pm0.000}$ | $0.627_{\pm0.059}$ |
| HEALNet (Inter.) | $0.714_{\pm0.025}$ | $0.618_{\pm0.063}$ | $0.842_{\pm0.063}$ | $0.594_{\pm0.025}$ | $0.767_{\pm0.022}$ | $0.748_{\pm0.009}$ | $0.639_{\pm0.09}$ |
| *LegoFuse, w/ 2 epochs* | $0.734_{\pm0.032}$ | $0.626_{\pm0.112}$ | $0.863_{\pm0.112}$ | $0.634_{\pm0.010}$ | $0.771_{\pm0.020}$ | $0.759_{\pm0.041}$ | $0.701_{\pm0.023}$ |

Mean and standard deviation of task performance, showing the concordance Index (survival) and AUC (classification) on 5 random sub-sampling folds with the best and second-best models highlighted.
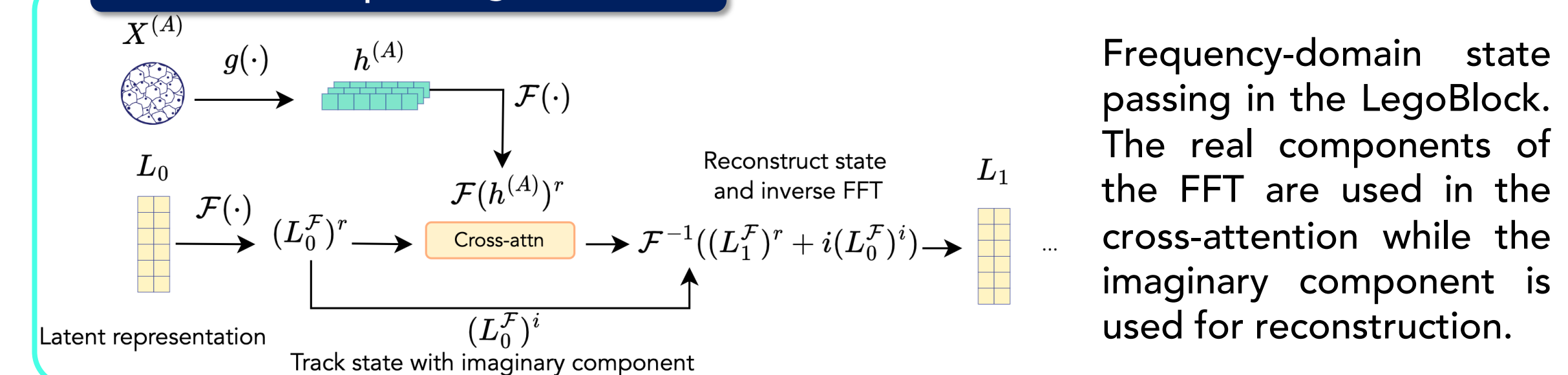
### Merge vs. Ensemble



AUC performance on the MIMIC dataset when merging pre-trained using **LegoMerge** and **LegoFuse**. Our multimodal model merge shows significantly improved performance over using and **ensemble**, exhibiting the performance gains at no additional cost through the merge.

### Latent state passing



Frequency-domain state passing in the LegoBlock. The real components of the FFT are used in the cross-attention while the imaginary component is used for reconstruction.

### LegoMerge outperforms unimodal models without additional training



BLCA — Uplift: 2.0% · BRCA — Uplift: 1.0% · KIRP — Uplift: -1.5% · UCEC — Uplift: 1.0% · ICD9 — Uplift: 3.9% · MORT — Uplift: 4.0% · ISIC — Uplift: 1.5%

Mean task performance (c-Index/AUC) of **LegoBlock (Tabular)**, **LegoBlock (Image/Time Series)** and **LegoMerge**, showing the increase in task performance by applying the multimodal model merge without any fine-tuning. Our proposed method shows improved performance on 6 out of 7 tasks.