

# VIRTUAL: Vector-based Dynamic Power Estimation via Decoupled Multi-Modality Learning

Yuntao Lu<sup>1</sup>, Mingjun Wang<sup>1,2</sup>, Yihan Wen<sup>1</sup>, Boyu Han<sup>3</sup>, Jianan Mu<sup>2</sup>, Huawei Li<sup>2</sup>, Bei Yu<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong

<sup>2</sup>Institute of Computing Technology, CAS

<sup>3</sup>Stanford University

October 27, 2025







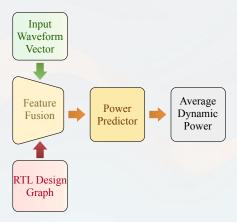
## Outline

- 1 At a Glance
- 2 Motivation
- 3 Problem Definition
- 4 Methodology
- **6** Evaluations
- 6 Conclusions & Future Work



### At a Glance

### Vector-based Dynamic Power Estimation via Decoupled Multi-Modality Learning



**Fig.** 1 An overview architecture of VIRTUAL.

#### **Innovations**

- Fast end-to-end power prediction.
  - Eliminate gate-level simulation.
  - Read input vector and RTL graph.
- Multi-modal decoupling and fusion.
  - Decouple input vectors & RTL graphs.
  - Eliminate toggle-rate computation.
- **Self-Supervised Contrastive Learning** 
  - Power-aware vector encoding.
  - Functionally-focused RTL encoding.

### Performance Highlights

| # RTL Designs | Pearson | MAPE   | Speedup |
|---------------|---------|--------|---------|
| 2004          | 0.84    | 23.43% | 14.27×  |



# Motivation: Traditional Power Analysis Workflow

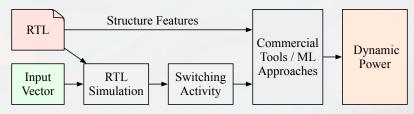


Fig. 2 Traditional power estimation workflow by commercial tools or ML approaches.

#### **Bottlenecks**

- Runtime: Simulation & synthesis hours to days.
- 2 Iteration: Full flow re-execution of each change.
- 6 Scalability: Larger designs, longer runtimes.

### **Impacts**

- Higher feedback latency.
- Limited optimization iterations.



# Motivation: Challenges & Limitations of Power Analysis

### Challenges of Traditional Power Analysis Flow

- RTL/Gate-level(Netlist) simulation.
- VCD/SAIF files for switching activity.
- 3 Toggle rate extraction from VCD/SAIF files.

### State-of-the-Art (SOTA) AI-driven Power Estimations

**Tab.** I SOTA AI-based approaches for power estimation.

| Method                 | Netlist Req. | Gate-Level Sim. | RTL Sim. | Toggle Rate Req. | Speed  |
|------------------------|--------------|-----------------|----------|------------------|--------|
| GRANNITE1              | Yes          | Yes             | No       | Yes              | Slow   |
| GRASPE <sup>2</sup>    | Partial      | Partial         | Partial  | Yes              | Medium |
| MasterRTL <sup>3</sup> | No           | No              | Yes      | Yes              | Medium |

⇒ Approaches rely on simulation results for modeling switching activities for each design.

<sup>3</sup>Wenji Fang et al. (2023). "MasterRTL: A Pre-synthesis PPA Estimation Framework for Any RTL Design". I Proc. ICCAD. IEEE, pp. 1–9.

<sup>&</sup>lt;sup>1</sup>Yanqing Zhang, Haoxing Ren, and Brucek Khailany (2020). "GRANNITE: Graph Neural Network Inference for Transferable Power Estimation". In: *Proc. DAC*. IEEE, pp. 1–6.

<sup>&</sup>lt;sup>2</sup>MB Rakesh et al. (2023). "GRASPE: Accurate Post-synthesis Power Estimation from RTL using Graph Representation Learning". In: *Proc. ISCAS*. IEEE, pp. 1–5.

# Motivation: Gaps for Fast and Accurate Power Estimation

#### Limitations of Current AI-driven Methods

- Weak input vector modeling (toggle rates, switching patterns).
- Require simulation for toggle rate extraction.
- Limited generalization across input vectors and RTL designs.

### **Bridging Gaps**

**Goal**: Enabling fast and accurate power estimation for diverse inputs and RTLs.

- Representations for input vectors and RTL designs.
- Direct RTL-level estimation without simulation.
- Generalize across diverse inputs and RTLs.

### **Problem Definition**

#### Problem Formula

#### Given

- **Input Waveform Matrix**:  $W \in [0, 1]^{M \times T}$ ,
  - *M*: # of input ports.
  - T: # of clock cycles.
- **Q** RTL Graph: G = (V, E).
  - V: Set of functional units.
  - E: Set of signal connections between units.

### Objective:

• Learn a mapping function  $F_{power}(W,G) \rightarrow \hat{P}$ 

### VIRTUAL: Framework Overview

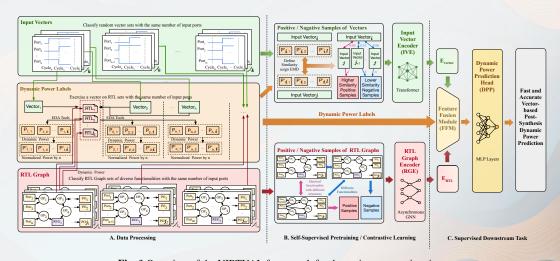


Fig. 3 Overview of the VIRTUAL framework for dynamic power estimation.

# VIRTUAL: Three-Stage Learning Strategy

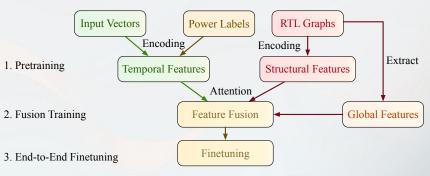


Fig. 4 Three-stage learning strategy of VIRTUAL.

### Stage 1: Pretraining

- Power-relevant features.
- Functional-focused features.

### Stage 2: Fusion Training

- Temporal features.
- Structural features.

### Stage 3: Fine-tuning

- Joint optimized modules.
- Supervised learning.



# Input Vector Encoding: Pattern Extraction

#### Transformer-Based Pattern Extraction

- **1** Tokenization:
  - Split fixed-length segments.
  - Add positional encoding.
- Multi-head Self-Attention:
  - Capture long-range dependencies.
  - Learn ports correlation.
  - Identify switching patterns.
- **6)** Feed-forward Networks:
  - Transform attention outputs.
  - Extract power-relevant temporal features.
- $\Rightarrow$  Temporal embedding  $E_{\text{vector}} \in \mathbb{R}^d$  capturing switching patterns.



# Input Vector Encoding: Power-Aware Contrastive Learning

### Vectors with similar power distributions should have similar embeddings

### **Pairing Strategy**

- Normalization: normalize power distributions for each vector by port count on multiple designs.
- Positive pairs: vectors with similar power distributions.
- Negative pairs: others.
- Learn to push apart dissimilar vectors in embedding views.

#### InfoNCE Loss Function

$$\mathcal{L}_{ ext{vector}} = -\sum \log rac{\exp(\cos(E_{ ext{anchor}}, E_{ ext{pos}})/ au)}{\sum \exp(\cos(E_{ ext{anchor}}, E_{ ext{neg}})/ au)}$$

, where  $\tau$  is temperature,  $\cos$  is cosine similarity.

#### **Benefits**

- Generalize across input waveforms.
- Capture switching patterns.
- Minimal labels required.



# RTL Graph Encoding: Structure Extraction

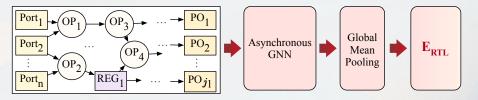


Fig. 5 GNN-based RTL functionality encoding.

### RTL Graph Representation

- Nodes: Simple functional units (registers, logic gates, operators).
- Edges: Signal connections (data flow).
- Node features: Unit type, fan-in/out, depth, bit-width, etc.



# RTL Graph Encoding: Asynchronous Message Passing

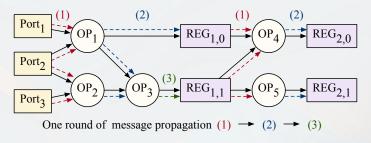


Fig. 6 Asynchronous message passing scheme mimicking signal propagation.

### Asynchronous GNN Architecture

- Message Passing: Mimic hardware signal propagation.
- Multi-head Graph Attention: Aggregate neighbor information.
- Temporal Updates: Model sequential signal flow.
- Global Pooling: Aggregate node embeddings to the graph-level.
- $\Rightarrow$  Structural embedding  $E_{RTL} \in \mathbb{R}^d$  capturing circuit functionality.



# RTL Graph Encoding: Functionality-Focused Contrastive Learning

### Functionally equivalent RTL designs should have similar embeddings.

### **Pairing Strategy**

- Positive pairs: Functionally equivalent designs
  - Same logic function, different optimizations.
  - Generated via Design Compiler with different settings.
- Negative pairs: Others.
- Learn to push apart dissimilar graphs in embedding views.

#### **Benefits**

- Generalize across RTL designs.
- Capture structural patterns.
- Learn without labels.

### VIRTUAL: Feature Fusion & Power Prediction

#### Feature Fusion Module

- Fusion Embeddings:
  - Temporal embedding:  $E_{\text{vector}} \in \mathbb{R}^d$
  - Structural embedding:  $E_{\text{graph}} \in \mathbb{R}^d$
  - Global features:  $E_{\text{global}} \in \mathbb{R}^k$
- Multi-head Cross-Attention:
  - Learn correlations of  $E_{\text{vector}}$  and  $E_{\text{graph}}$ .
  - Weighted combination of both modalities.
- **6)** Concatenation:
  - Combine with global features.
  - Unified representation:  $E_{\text{fused}} \in \mathbb{R}^{2d+k}$

#### Inference

• End-to-End Prediction: RTL graph G + Input vectors  $W \rightarrow$  average dynamic power.

#### **Power Prediction Head**

- **1) Architecture**: 2-layer fully-connected neural network.
- Supervised Training:

$$\mathcal{L}_{\text{power}} = \text{MSE}(\hat{P}, P) = \frac{1}{N} \sum_{i=1}^{N} (\hat{P}_i - P_i)^2$$

**6** End-to-End Fine-tuning:

$$\mathcal{L}_{total} = \mathcal{L}_{power} + \lambda_1 \mathcal{L}_{vector} + \lambda_2 \mathcal{L}_{graph}$$



# **Evaluation: Experimental Setup**

#### **Datasets**

- Circuits: 2004 RTL designs from OpenCores, HuggingFade, and Github.
  - Scalability:  $\sim 500$  to  $\sim 30,000$  nodes.
  - Functionality: arithmetics, controllers, communication, etc.
- **Input Vectors**: 100 vectors per circuit.
  - $\sim$  60 to  $\sim$  1,000 input ports.
  - 1,000 clock cycles per vector.
- **Ground Truth**: Average dynamic power values (switching + internal) of Design Compiler.

#### Metrics

- PCC: Correlation measurement.
- $\mathbf{R}^2$ : Goodness of fitness.
- MAPE: Relative error.
- **Speedup**: Runtime over traditional flow.



# **Evaluation: Performance Summary**

**Tab.** II Performance comparison between VIRTUAL AND SOTA AI-based methods

| Method                  | Circuit-based |       |        | Vector-based |       |        |
|-------------------------|---------------|-------|--------|--------------|-------|--------|
| Method                  | PCC           | $R^2$ | MAPE   | PCC          | $R^2$ | MAPE   |
| MasterRTL               | 0.786         | 0.752 | 30.85% | 0.757        | 0.692 | 30.79% |
| GRASPE                  | 0.815         | 0.775 | 29.72% | 0.780        | 0.647 | 30.41% |
| VIRTUAL w/o Contrastive | 0.757         | 0.729 | 29.09% | 0.745        | 0.765 | 26.48% |
| VIRTUAL (Full)          | 0.842         | 0.798 | 23.43% | 0.834        | 0.776 | 24.56% |

- Strong correlation with ground truth (PCC: 0.842).
- Relative low error across designs (MAPE: 23.43%).
- Significant speedup over the traditional workflow  $(14\times)$ .

#### Observations

- VIRTUAL outperforms SOTA AI-based methods in PCC and MAPE.
- Contrastive learning improves performance.
- Strong generalization across designs and inputs.

### Conclusions & Future Work

#### Conclusions

- **Runtime**: Eliminates synthesis/simulation for power analysis.
- Effectiveness: Demonstrates the effectiveness of multi-modal contrastive learning.
- Rapid Feedback: Enables design optimizations with rapid feedback.

#### Future Work

- **Scalability**: Handle larger designs (thousands of nodes).
- Generalization: Few-shot learning for unseen design families.
- Transferability: Transfer learning across different technology nodes.

# Thanks!

Q&A

