

1. Outline

Multimodal Learning for Recognition

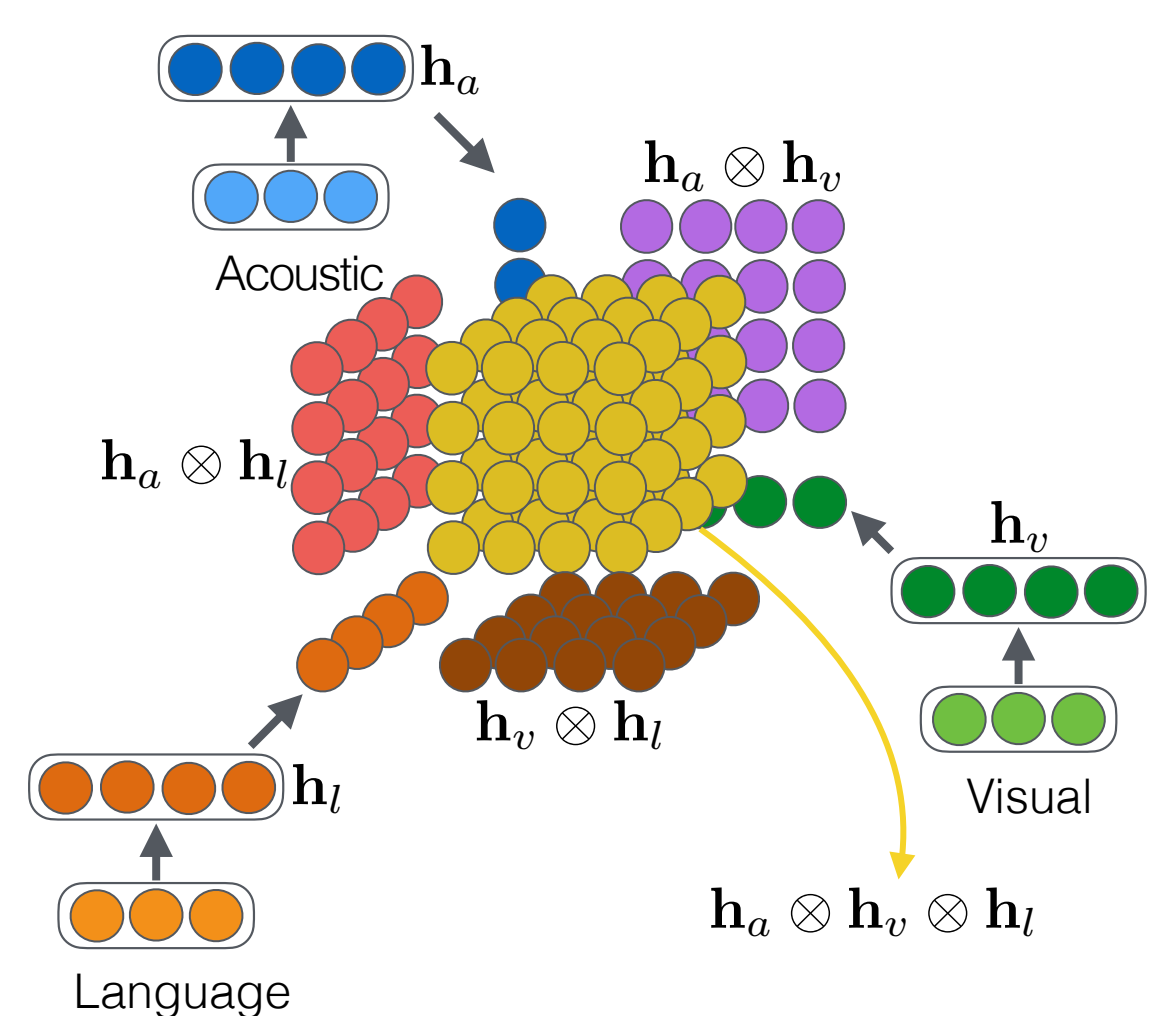
- Multimodal recognition integrates features of multiple modalities (language, acoustic and visual) for yielding robust predictions.
- Multimodal fusion** is a key step of multimodal recognition.
- Tensor-based fusion methods have achieved a great success.

Limitations of existing tensor fusion

- Restrict interaction among modalities to be **linear w.r.t. each modality**
- Ignore **high-order statistical information**

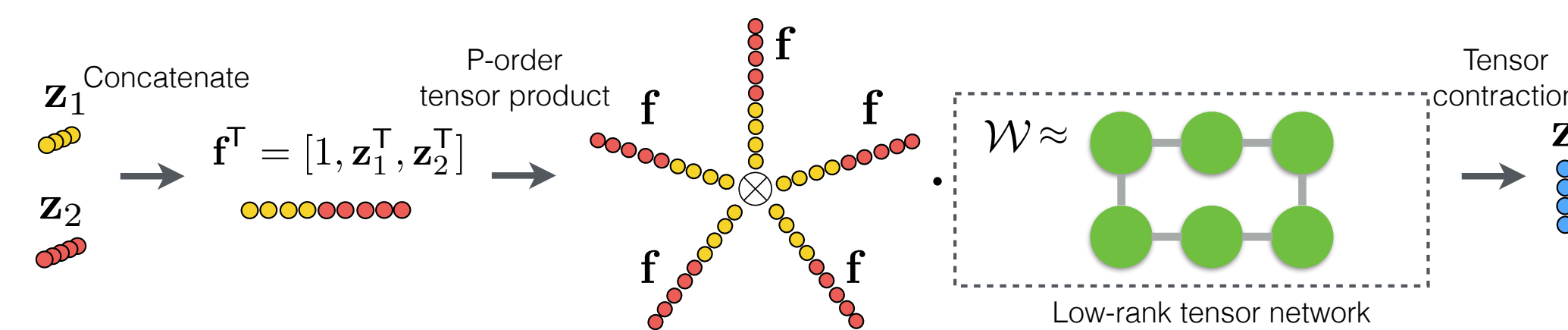
Our contributions

- Explicitly model nonlinear intra-modal and cross-modal interactions via **high-order polynomial moments**
- Directly model **local interactions across mixed dimensions over time**
- Significantly reduce heavy computation via using tensor network



2. Polynomial Tensor Pooling (PTP)

- PTP** block first fuses M feature vectors using high-order moments and then transforms them into a joint representation.



- Fusion via high-order moment is obtained by P -order tensor product of concatenated features: $f^T = [1, z_1^T, z_2^T, \dots, z_M^T]$

$$\mathcal{F} = \underbrace{\mathbf{f} \otimes \mathbf{f} \otimes \dots \otimes \mathbf{f}}_{P\text{-order}}$$

- Transformation is performed by a weight tensor: $\mathcal{W} = [\mathcal{W}^1, \dots, \mathcal{W}^h, \dots, \mathcal{W}^H]$

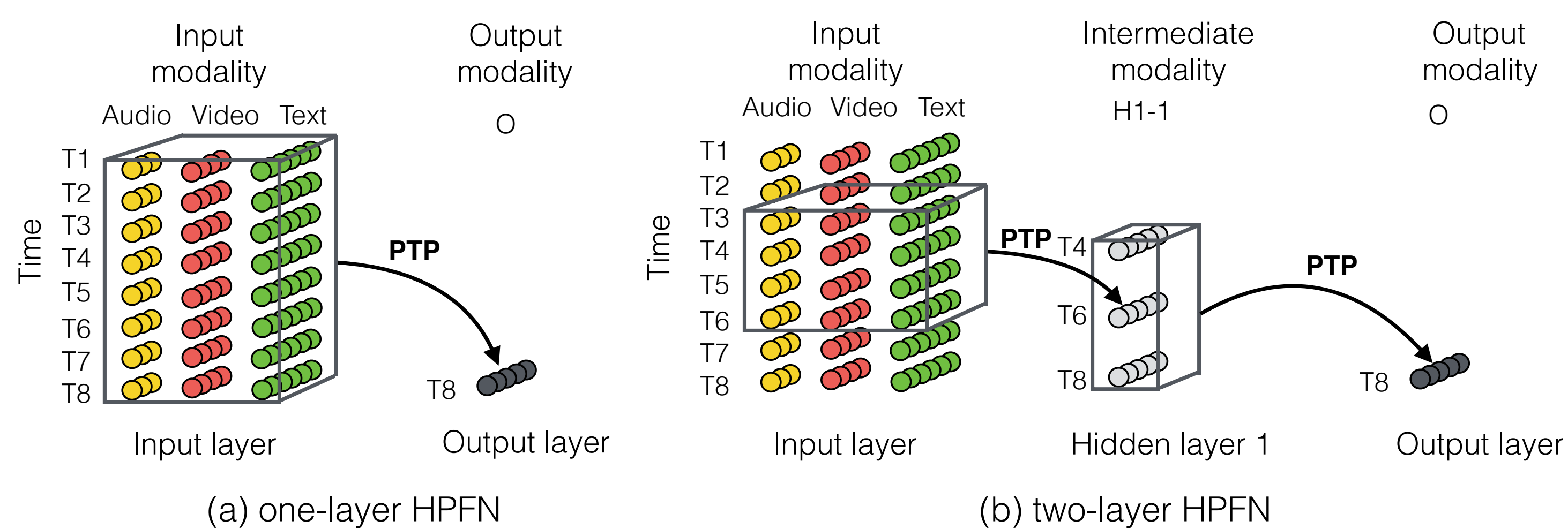
$$z_h = \sum_{i_1, i_2, \dots, i_P} \mathcal{W}_{i_1 i_2 \dots i_P}^h \cdot \mathcal{F}_{i_1 i_2 \dots i_P}$$

- Low-rank **tensor networks** is used to reduce large computation.

$$z_h = \sum_{i_1, i_2, \dots, i_P} \left[\left(\sum_{r=1}^R a_r^h \prod_{p=1}^P w_{r i_p}^h \right) \left(\prod_{p=1}^P \mathbf{f}_{i_p} \right) \right] = \sum_{r=1}^R a_r^h \prod_{p=1}^P \sum_{i_p} w_{r i_p}^h \mathbf{f}_{i_p}$$

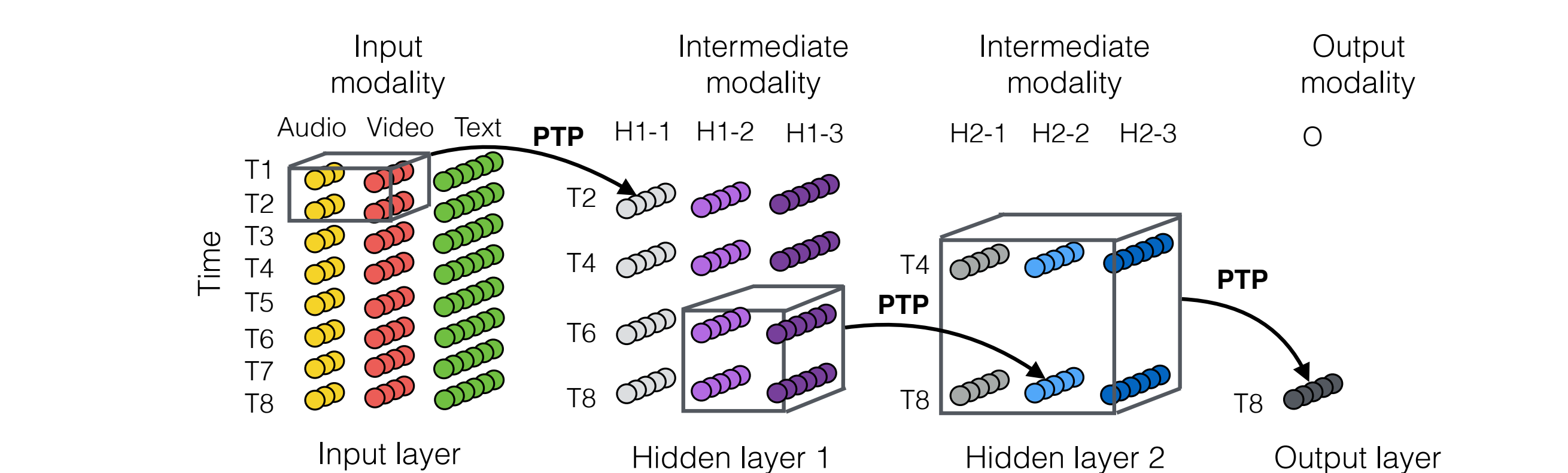
3. Hierarchical Polynomial Fusion Network (HPFN)

- Features of multiple modalities are rearranged into a **“feature map”**.
- Single-layer HPFN** is constructed by a global PTP operating on a **“receptive window”** across all time steps and modalities.

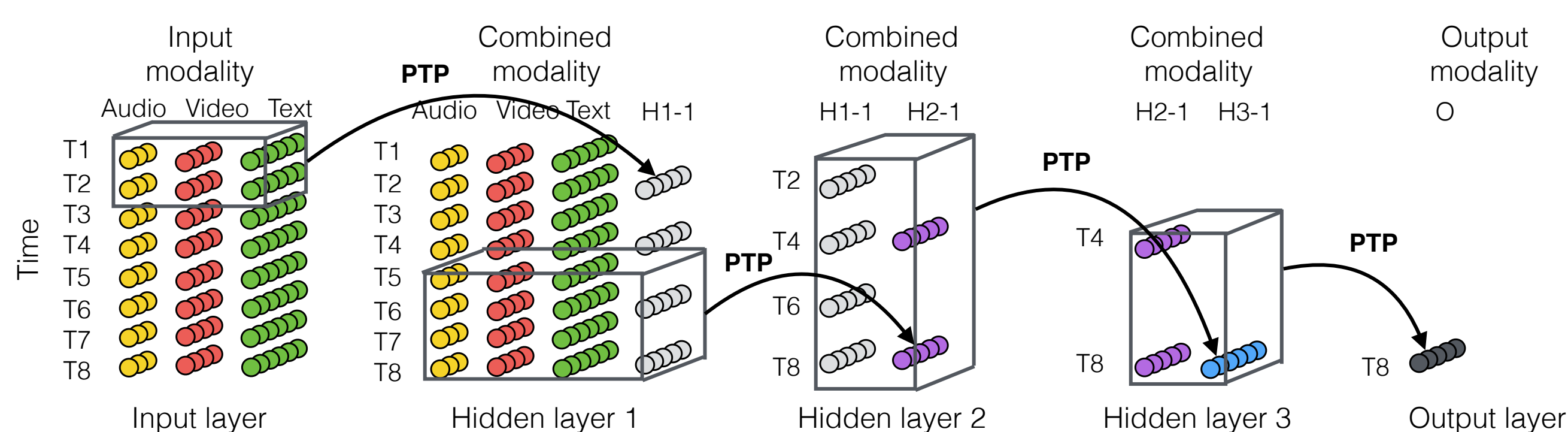


- Multi-layer HPFN** is established by recursing PTP blocks layer by layer into a tree-structured architecture.

- Local temporal-modality interactions most relevant to prediction can be transmitted to the global level.
- PTP can be treated as a **“fusion filter”** analogous to a CNN filter.
- CNN-style fusion framework** with flexible design choices bring benefits.



- Incorporate **dense connectivity** to enhance the expressive capacity

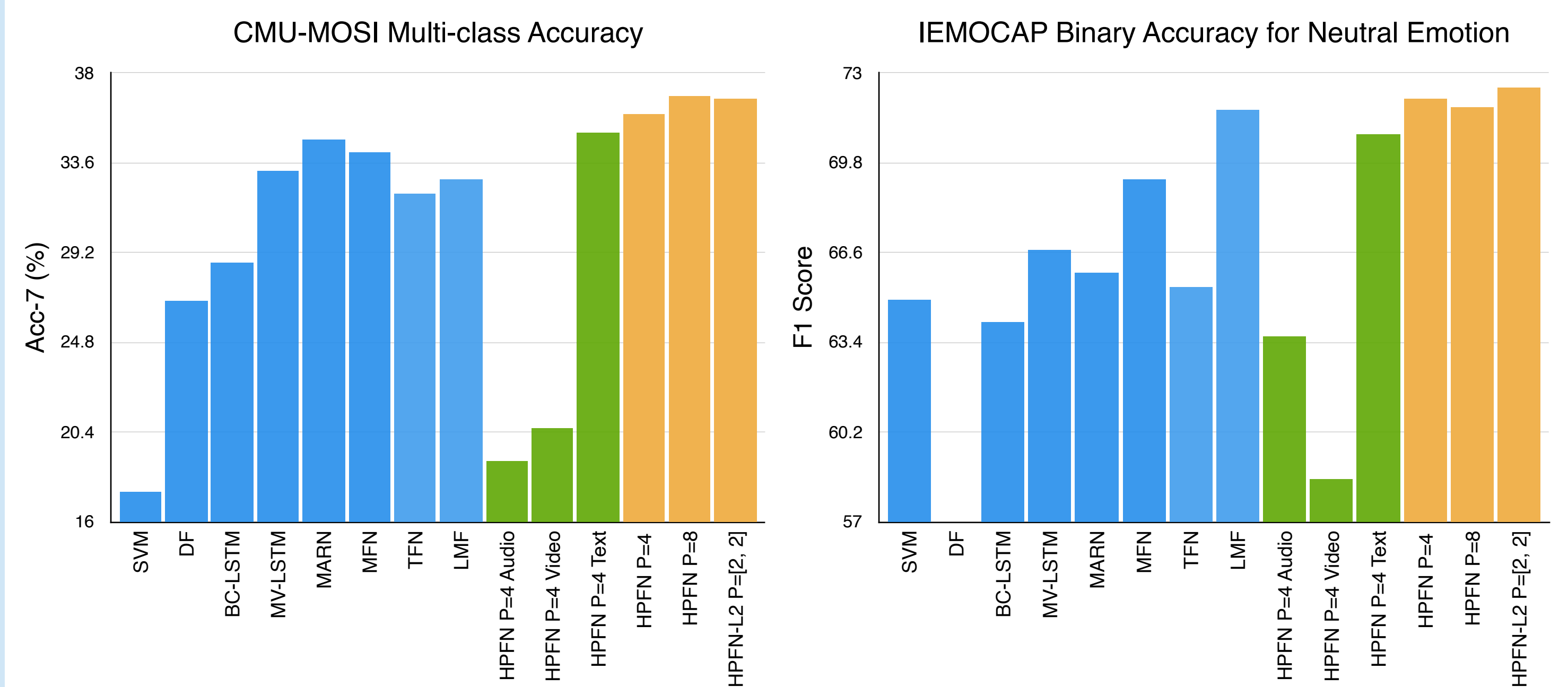


- Model complexity depends on total number of PTP filters and window size, parameters are larger than **LMF** but much smaller than **TFN**.

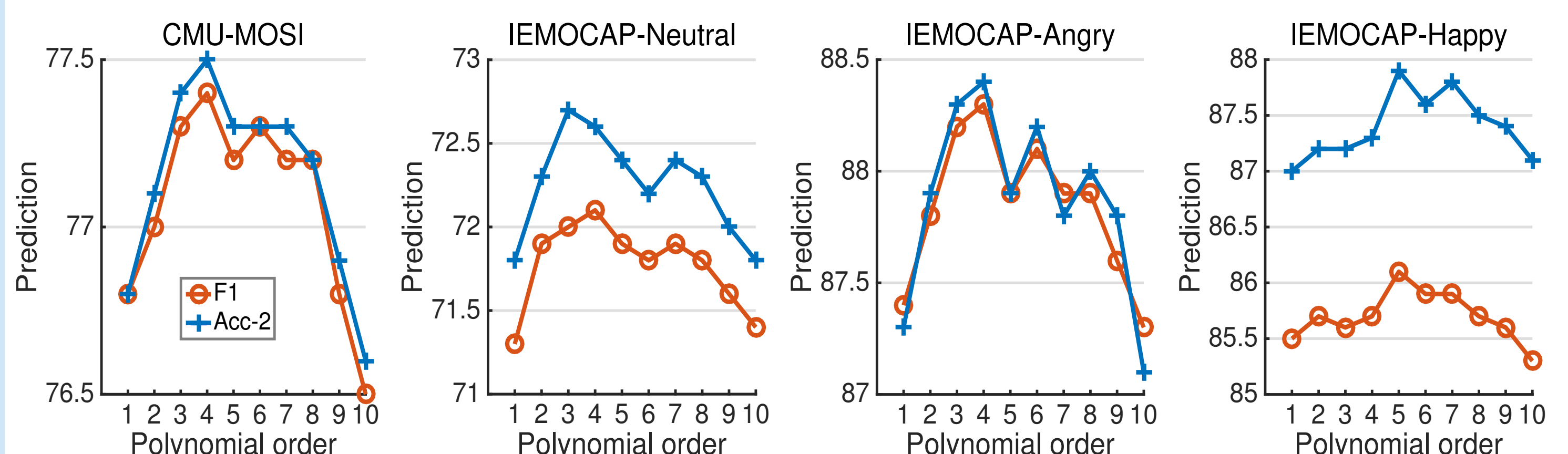
Model	TFN [non-temporal]	LMF [non-temporal]	PTP [temporal]	HPFN (L layers) [temporal]
Param.	$\mathcal{O}(I_y \prod_{m=1}^M I_m)$	$\mathcal{O}(I_y R (\sum_{m=1}^M I_m))$	$\mathcal{O}(I_y R (\sum_{t=1}^T \sum_{m=1}^M I_{t,m}))$	$\mathcal{O}(I_y R (\sum_{l=1}^L N_l) (\sum_{t=1}^T \sum_{m=1}^M I_{t,m}))$

4. Experimental Results

- Datasets and tasks:
 - CMU-MOSI** utterance level multimodal sentiment analysis with intensity range in $[-3, 3]$
 - IEMOCAP** utterance level binary classification for emotions, including neutral, angry, happy and sad
- Accuracy comparisons of HPFN with other methods



- Effect of fusion of one-layer HPFN with higher order polynomial



- Effect of fusion of HPFN with depth, connectivity

- Effect of fusion with and without the incorporation of temporal factors

