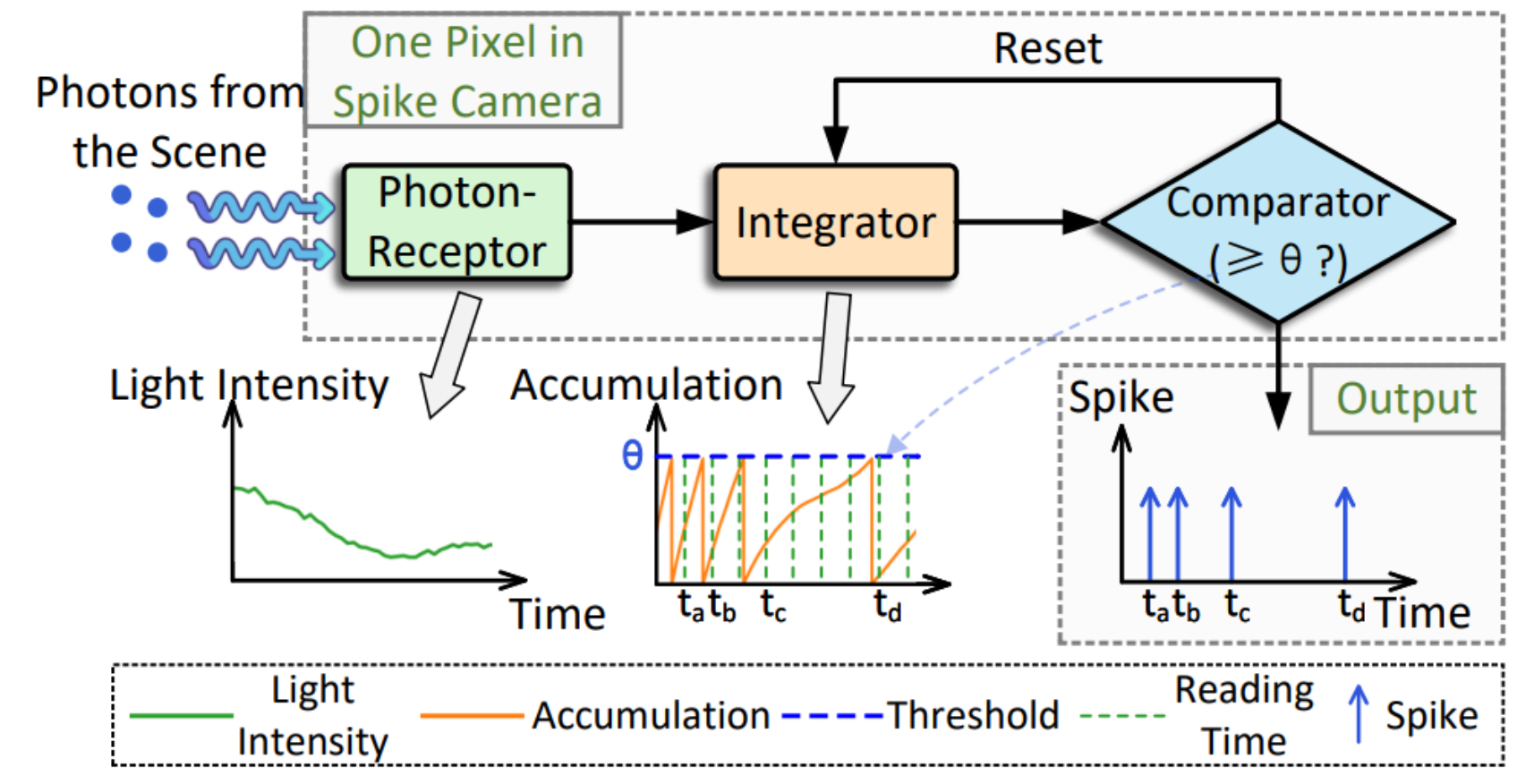


1. Introduction

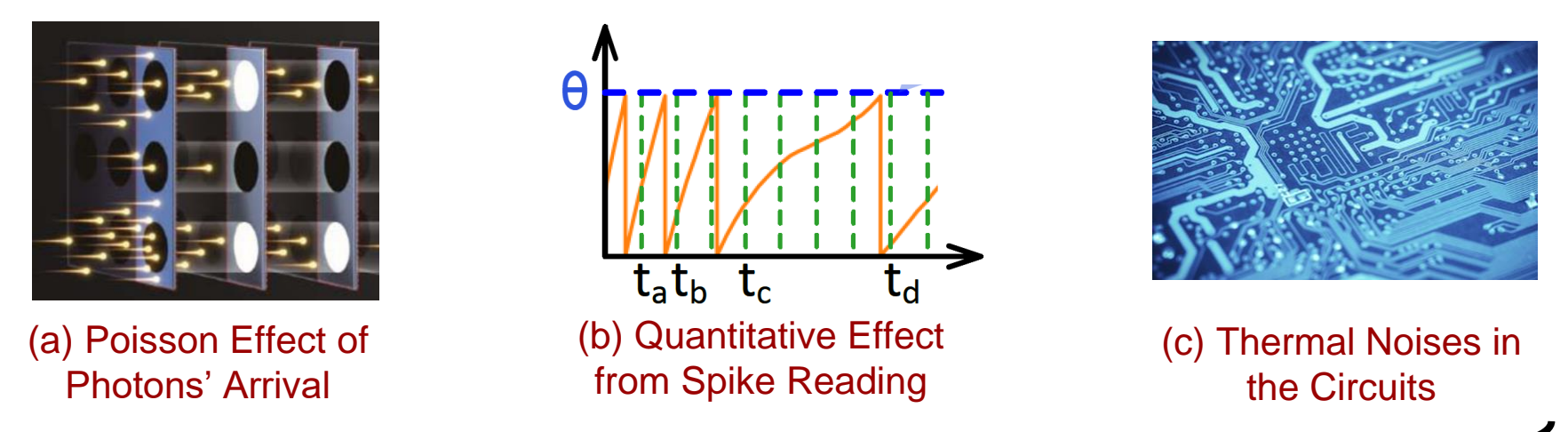
1.1 Spike Camera. Spike cameras are composed of an array of pixels working asynchronously. Each pixel of a spike camera is composed of three main components: photon-receptor, integrator, and comparator. The integrator accumulates the photoelectrons from the photon-receptor and transfers them to the voltage. The comparator compares the accumulation with the threshold continuously. Once the voltage of the integrator exceeds a certain threshold, the camera fires a spike and resets the accumulation.



$$A(\mathbf{x}, t) = \int_0^t \alpha \cdot I(\mathbf{x}, \tau) d\tau \bmod \theta$$

1.2 Challenges of Spike-Based Optical Flow.

Noises in the imaging of spike cameras.



Fluctuations and Randomness in Spikes

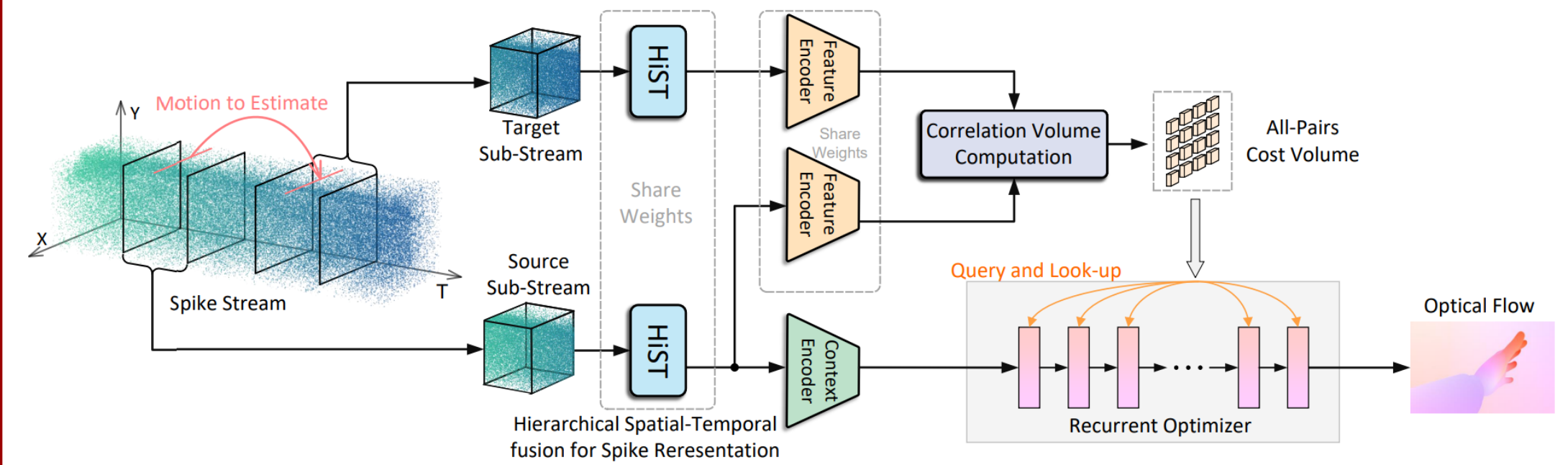
Ambiguities in correlation → Inaccurate feature matching

2. Contributions

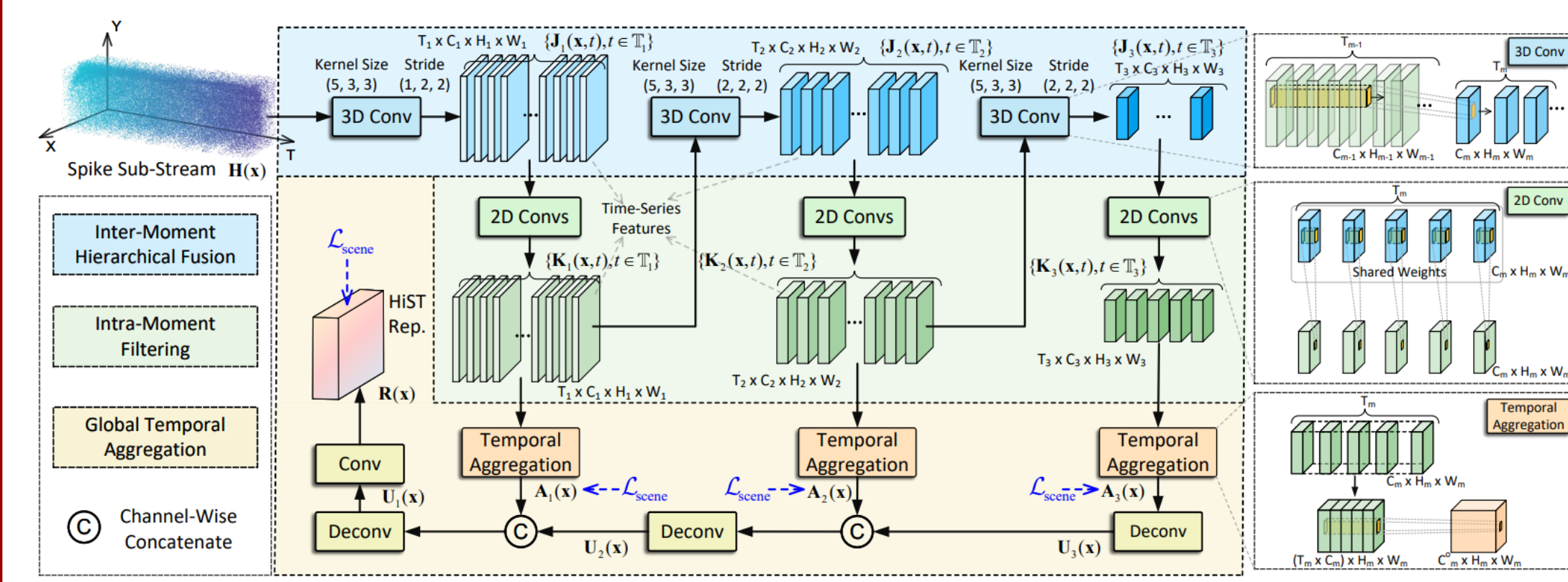
- A HiST-SFlow is proposed for spike-based optical flow. In HiST-SFlow, the spikes are represented by the HiST module and extracted to features for correlation. The optical flow is estimated by a recurrent optimizer.
- An inter-moment hierarchical fusion (InterF) module and an intra-moment filtering (IntraF) module are proposed to suppress the randomness in the spikes. A scene loss is proposed to constrain high-fidelity representation to contain the brightness information of the scene.

3. Approaches

3.1 Overall Architecture of HiST-Sflow.



3.2 Hierarchical Spatial-Temporal (HiST) Fusion.



(a) Inter-Moment Hierarchical Fusion (InterF).

- Fuse features at different moments while retaining the time information in features.

$$\mathbf{J}_m(\mathbf{x}, t) = \mathcal{J}_m[\{\mathbf{K}_{m-1}(\mathbf{x}, \tau) \mid \tau \in \mathbb{T}_{m-1}\}]$$

$$\mathbb{T}_{m-1} = \{T_c - T_{m-1}^{\text{half}}, \dots, T_c, \dots, T_c + T_{m-1}^{\text{half}}\}$$

(b) Intra-Moment Filtering (IntraF).

- Reduce the influence of spikes' fluctuations for each moment through the feature at the current moment.
- The InterF and IntraF are implemented alternatively in each level of the pyramid.

$$\mathbf{K}_m(\mathbf{x}, t) = \mathcal{K}_m[\mathbf{J}_m(\mathbf{x}, t)], t \in \mathbb{T}_m$$

(c) Global Temporal Aggregation (GTA).

- Fuse features of all the moments at each level of the pyramid to represent the central moment of input spike sub-stream.

$$\mathbf{A}_m(\mathbf{x}) = \mathcal{A}_m[\text{Cat}\{\mathbf{K}_m(\mathbf{x}, \tau) \mid \tau \in \mathbb{T}_m\}]$$

(d) Scene Loss.

- Ensure the spike representation contain the scene's brightness information
- The $\{\mathcal{P}_m\}_{m=0}^3$ are 3-layer convolution layers, which are used only during training and not in inference.

$$\mathcal{L}_{\text{scene}} = \|\mathbf{I}_{\text{scene}}(\mathbf{x}, T_c) - \mathcal{P}_0(\mathbf{R}_{T_c}(\mathbf{x}))\|_1 + \sum_{m=1}^3 \lambda_m \|\sigma_m(\mathbf{I}_{\text{scene}}(\mathbf{x}, T_c)) - \mathcal{P}_m(\mathbf{A}_m(\mathbf{x}))\|_1$$

3.3 Loss Function.

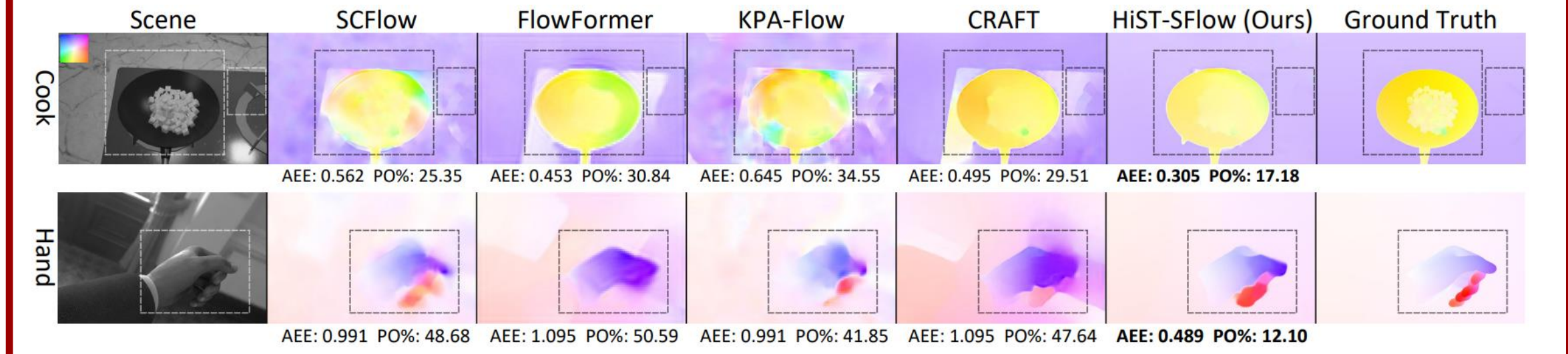
$$\mathcal{L}_{\text{flow}} = \sum_{i=1}^N \gamma^{N-i} \|\mathbf{w}_i(\mathbf{x}) - \mathbf{w}_{\text{gt}}(\mathbf{x})\|_1 \quad \mathcal{L} = \mathcal{L}_{\text{flow}} + \lambda(\mathcal{L}_{\text{scene}}^{\text{src}} + \mathcal{L}_{\text{scene}}^{\text{tgt}})$$

4. Experimental Results

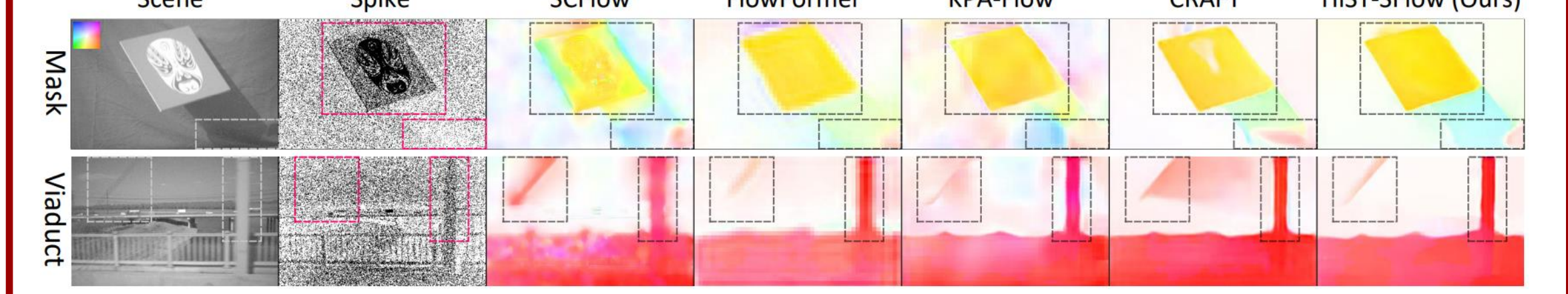
4.1 Quantitative Results on PHM Dataset (AEPE / PO%).

Architecture	Ball	Cook	Dice	Doll	Fan	Hand	Jump	Poker	Top	Average
SCFlow	0.51 / 20.3	1.34 / 38.6	1.10 / 30.7	0.22 / 5.6	0.24 / 10.7	1.30 / 57.3	0.11 / 3.0	0.80 / 41.1	2.14 / 17.7	0.863 / 25.00
RAFT	0.46 / 12.5	1.32 / 43.7	0.95 / 29.3	0.24 / 6.7	0.28 / 12.7	1.11 / 45.1	0.11 / 3.0	0.67 / 37.1	2.19 / 19.7	0.813 / 23.30
GMA	0.61 / 21.7	1.84 / 74.7	1.13 / 34.2	0.39 / 9.4	0.36 / 12.1	2.13 / 80.6	0.17 / 2.8	0.88 / 43.5	2.29 / 23.6	1.087 / 33.63
Flow1D	0.79 / 51.4	1.28 / 50.8	1.15 / 47.9	0.27 / 6.3	0.28 / 11.0	1.86 / 83.1	0.13 / 3.4	0.85 / 50.1	2.19 / 17.7	0.979 / 35.76
KPA-Flow	0.47 / 14.9	1.41 / 45.9	0.87 / 29.9	0.27 / 7.1	0.29 / 12.7	1.19 / 47.7	0.12 / 3.0	0.65 / 36.6	2.19 / 19.4	0.827 / 24.12
GMFlow	0.76 / 42.4	1.29 / 61.0	1.54 / 81.7	0.31 / 8.4	0.43 / 14.1	1.83 / 65.0	0.30 / 3.7	0.95 / 54.2	2.29 / 23.3	1.077 / 39.33
GMFlowNet	0.45 / 12.1	1.22 / 43.8	1.02 / 32.9	0.35 / 7.8	0.25 / 10.7	1.53 / 65.3	0.12 / 3.2	0.65 / 31.5	2.18 / 17.5	0.863 / 24.98
CRAFT	0.61 / 15.0	1.28 / 43.5	0.93 / 27.6	0.19 / 5.0	0.25 / 10.2	1.67 / 73.3	0.10 / 2.6	0.56 / 23.1	2.15 / 15.1	0.860 / 23.94
FlowFormer	0.52 / 13.5	1.48 / 58.7	0.98 / 31.0	0.25 / 6.7	0.29 / 11.5	1.82 / 84.5	0.14 / 3.6	0.94 / 54.9	2.22 / 19.5	0.959 / 31.54
HiST-SFlow	0.28 / 7.8	0.80 / 27.4	0.85 / 23.3	0.20 / 5.6	0.27 / 12.8	0.64 / 21.7	0.08 / 2.5	0.53 / 23.9	2.11 / 14.8	0.640 / 15.54

4.2 Visualization Results on PHM Dataset.



4.3 Visualization Results on Real-Captured Data.



4.4 Ablation Studies.

Index	Ablations on Proposed Modules				Ablations on Different Representations				
	Settings	$\Delta t = 10$		$\Delta t = 20$		$\Delta t = 10$		$\Delta t = 20$	
(A)	InterF	AEPE	PO%	AEPE	PO%	AEPE	PO%	AEPE	PO%
(B)	IntraF	0.986	33.17	2.095	56.56	0.868	25.72	1.757	34.19
(C)	$\mathcal{L}_{\text{scene}}$	0.694	18.18	1.449	21.99	0.880	29.77	1.824	37.91
(D)	InterF	0.676	17.34	1.433	22.79	0.799	21.10	1.703	34.58
(E)	IntraF	0.675	16.63	1.448	21.40	0.696	16.99	1.533	23.36
(F)	GTA	0.640	15.54	1.417	19.73	0.640	15.54	1.417	19.73

4.5 Using HiST for Other Baselines.

