

# Zero-Shot Event-Intensity Asymmetric Stereo via Visual Prompting from Image Domain

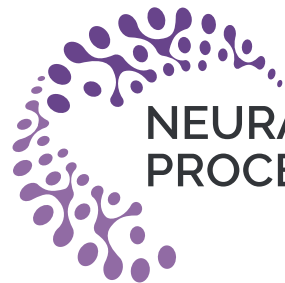
Hanyue Lou<sup>#1</sup> Jinxiu Liang<sup>#1</sup> Mingguai Teng<sup>1</sup> Bin Fan<sup>1</sup> Yong Xu<sup>2</sup> Boxin Shi<sup>1</sup>

<sup>1</sup> Peking University      <sup>2</sup> South China University of Technology

{hylz, cssherryliang, minggui\_teng, binfan, shiboxin}@pku.edu.cn yxu@scut.edu.cn



北京大學  
PEKING UNIVERSITY



NEURAL INFORMATION  
PROCESSING SYSTEMS



華南理工大學  
South China University of Technology

# Stereo vision



- Stereo vision estimates depth by mimicking human binocular vision.
- It computes the disparity between images captured by each camera to estimate the distance of objects.
- Applications: 3D reconstruction, robotics, autonomous driving.....

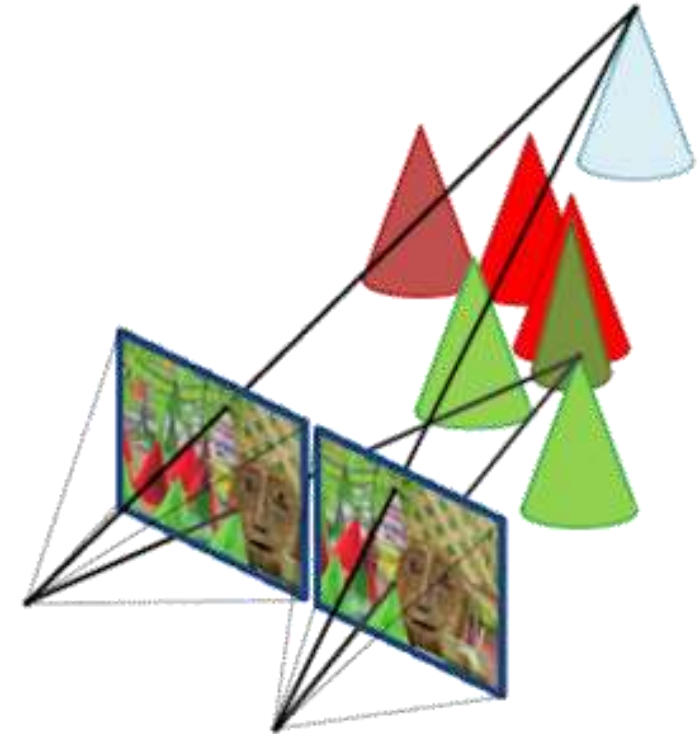


Image source: [en.ids-imaging.com](http://en.ids-imaging.com)



# Event-based vision



- Event cameras are novel visual sensors that report intensity changes, providing high dynamic range and high temporal resolution.
- However, no event signals are triggered when the scene is static or lacks texture.



Intensity images & Event signals

\*Video courtesy of Elias Mueggler

# Event-Image Stereo



- Image-based cameras suffer from low dynamic range and low temporal resolution. However, they always provide spatially dense information.
- Events and images provide complementary information, making them a good combination for a stereo system.

(Left view)



+



(Right view)

# Challenge



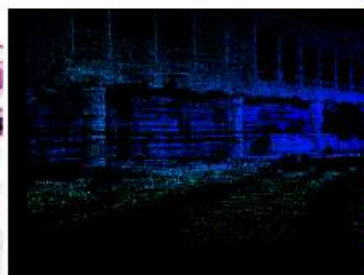
- Problem: Existing E.-I. A. S. datasets are not sufficient.
- Traditional methods → Limited performance
- Data-driven methods → Overfit on DSEC / MVSEC



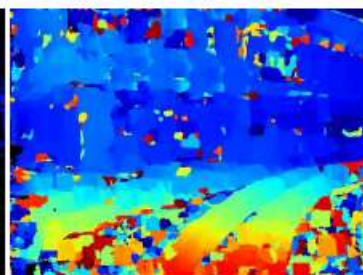
Frame (Left)



Event (Right)

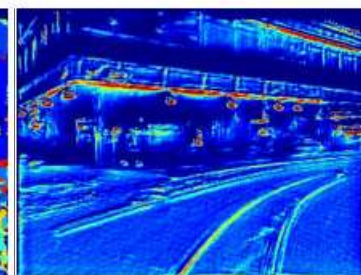


Ground Truth



HSM

↑  
Traditional



DAEI-MVSEC

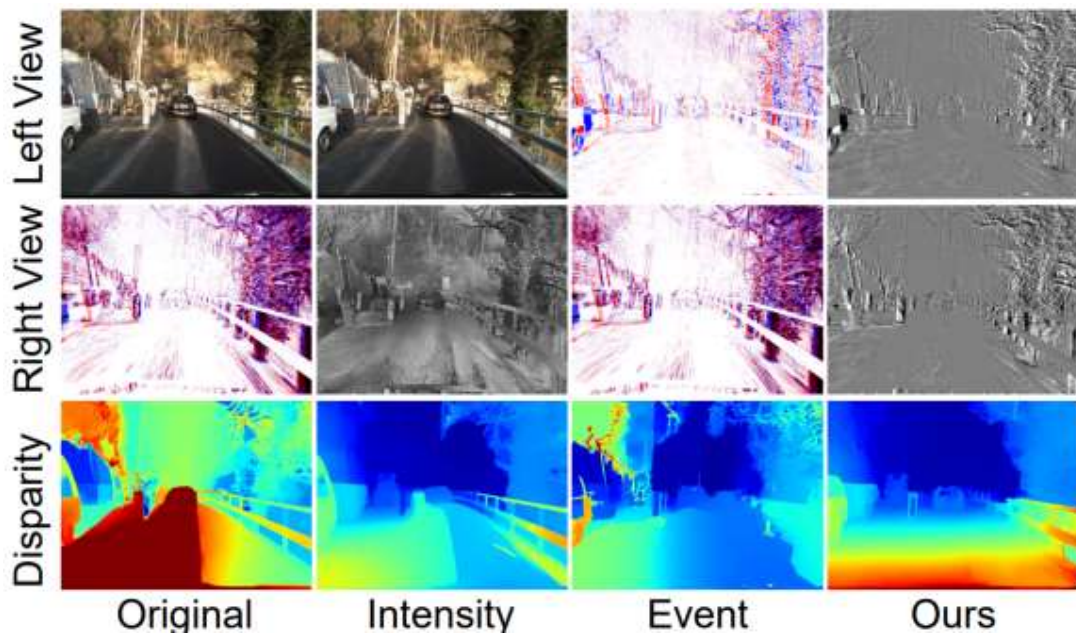
↑  
Data-driven



# Motivation



- Surprisingly: Existing frame-based methods work **very** well!
- Frame-based SOTA methods can generalize to non-natural images and even overcome large differences in appearance.

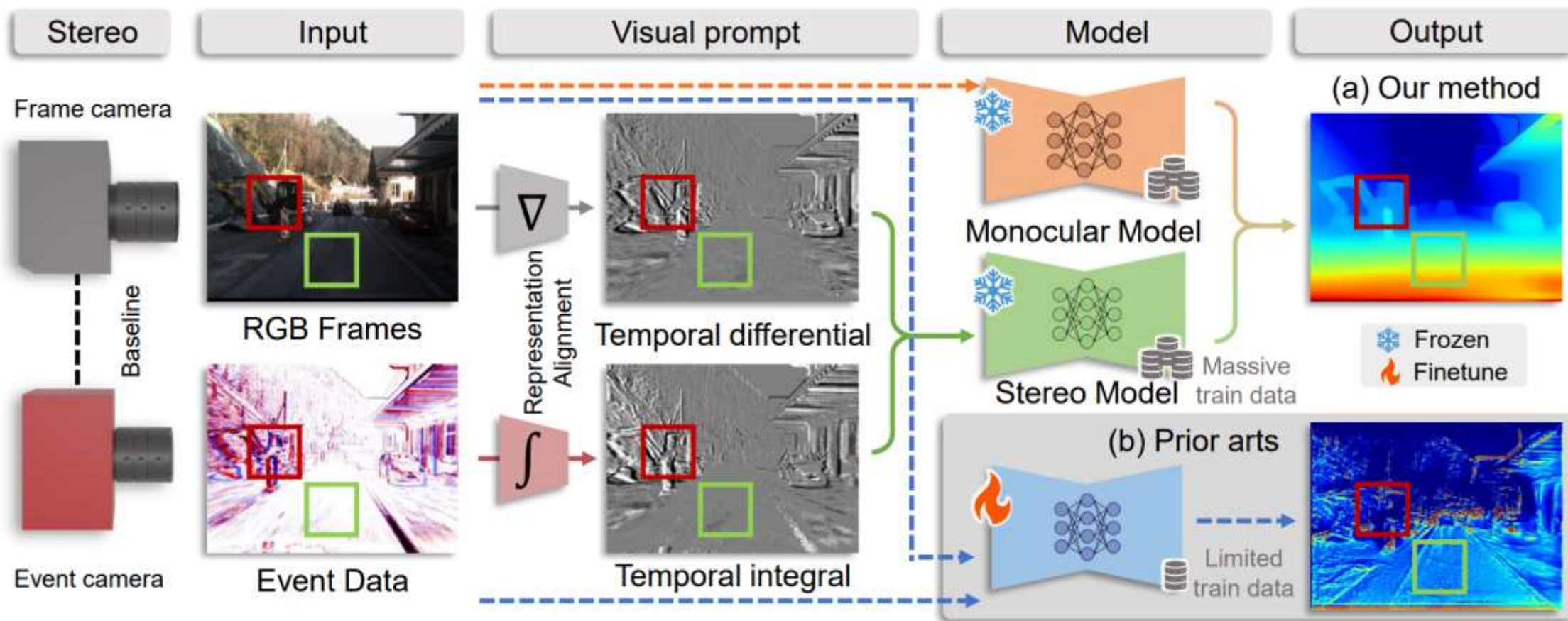


Frame-based stereo methods



E-I asymmetric stereo methods

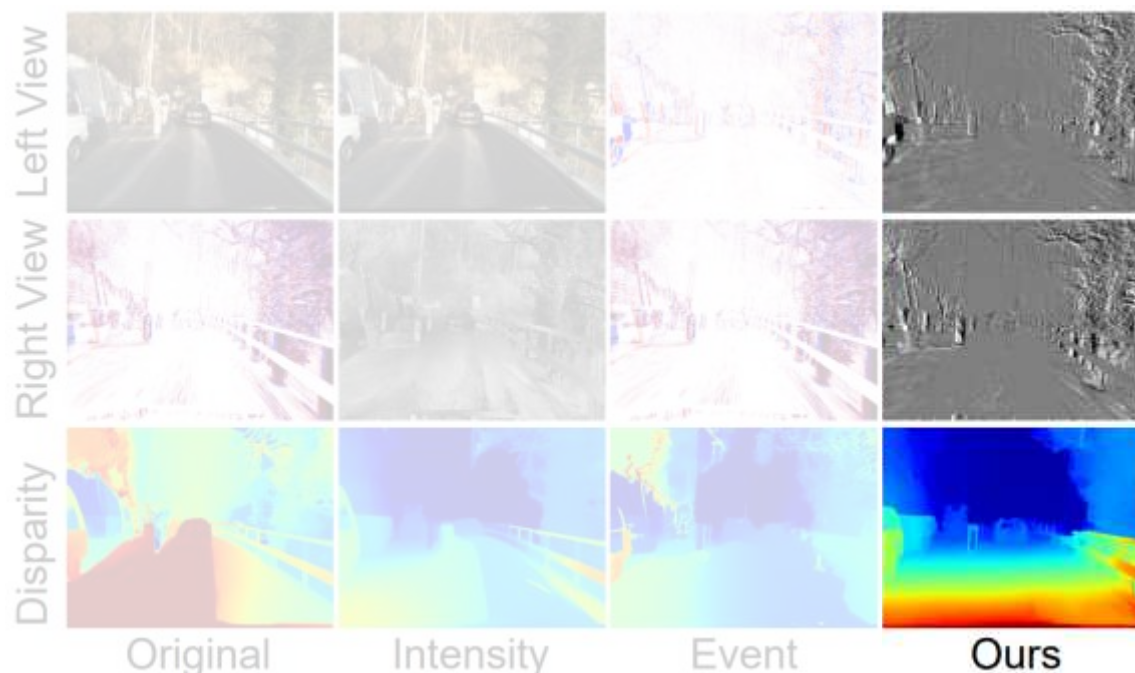
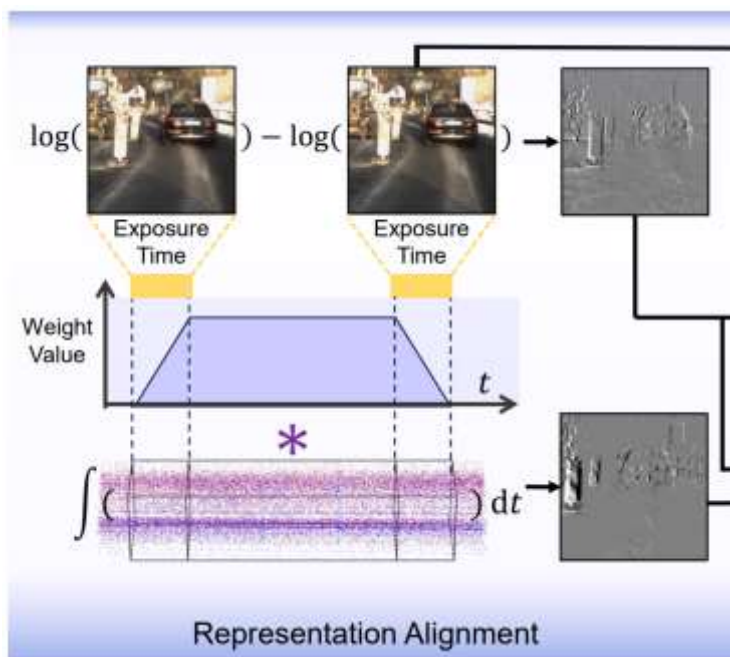
# Method



# Method



- We design a “visual prompt”: an intermediate representation that minimizes the appearance gap between images and event streams.

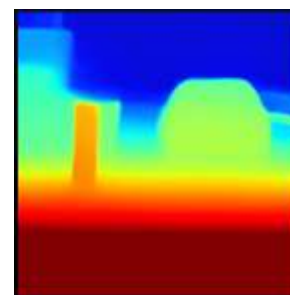
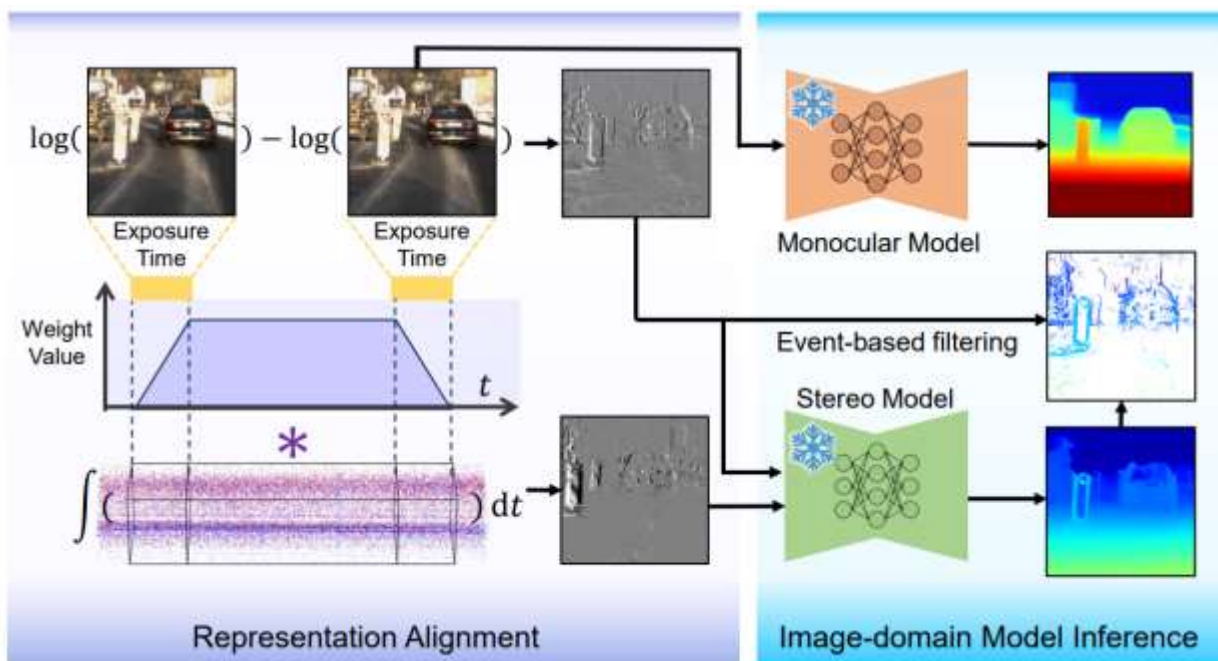




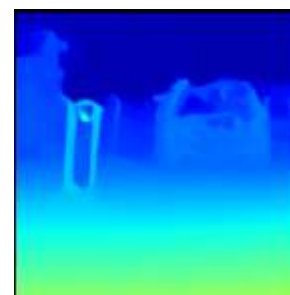
# Method



- However, the stereo method is erroneous where events are sparse. A frame-based monocular depth estimation model can provide complementary information, but its output is relative.



Monocular results  
😊 Dense  
😞 Relative depth

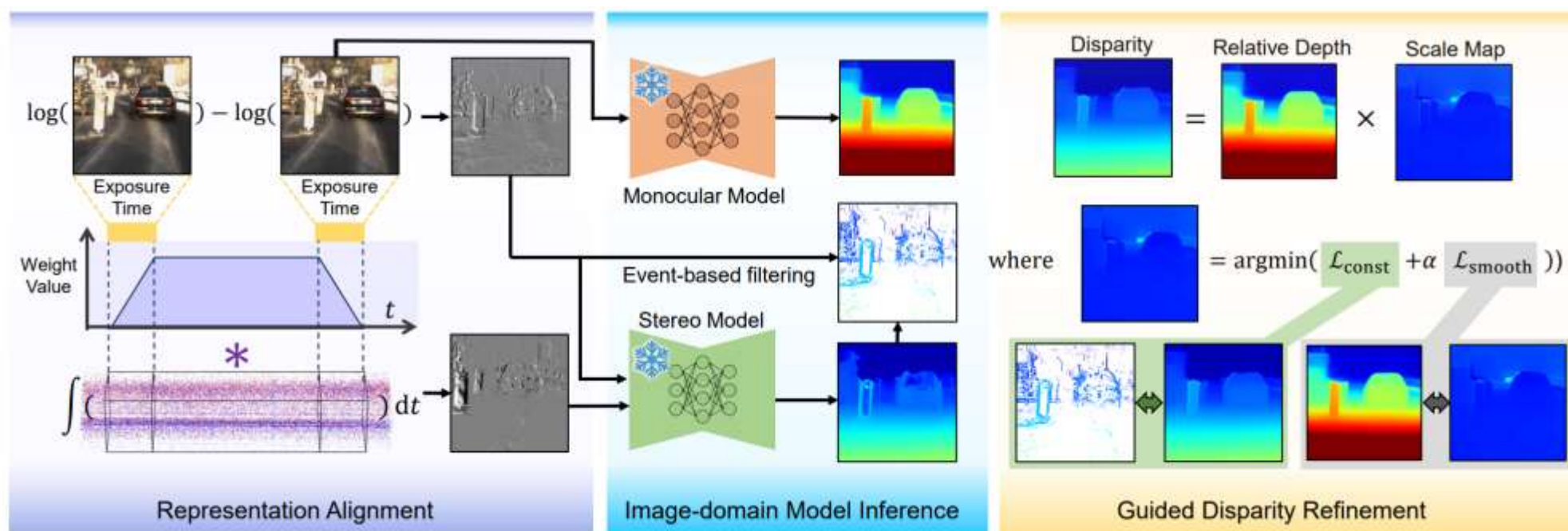


Monocular results  
😊 Absolute depth  
😞 Sparse

# Method



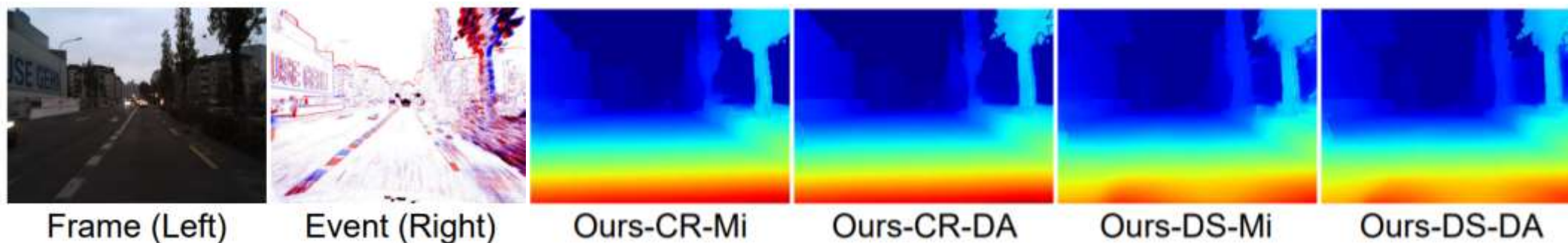
- We solve an optimization problem to fuse the relative monocular results and the absolute stereo results. The refined disparity only follows the stereo where events are dense.



# Results



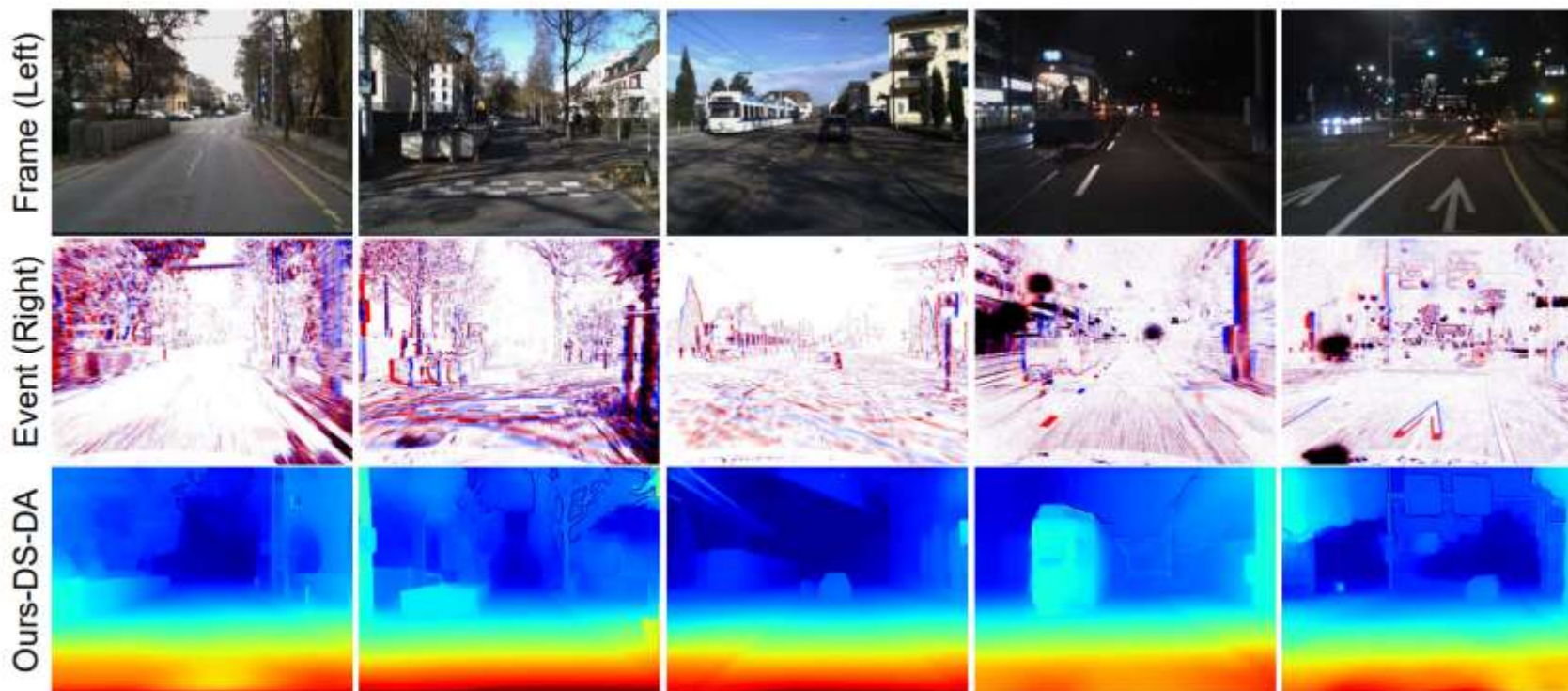
- Our framework achieves state-of-the-art performance among all zero-shot solutions.
- The frame-based stereo and monocular models used in the framework can be seamlessly changed without any finetuning, allowing for flexible upgrades as related fields advance.



# Results



- Our method performs robustly in diverse scenarios and datasets.

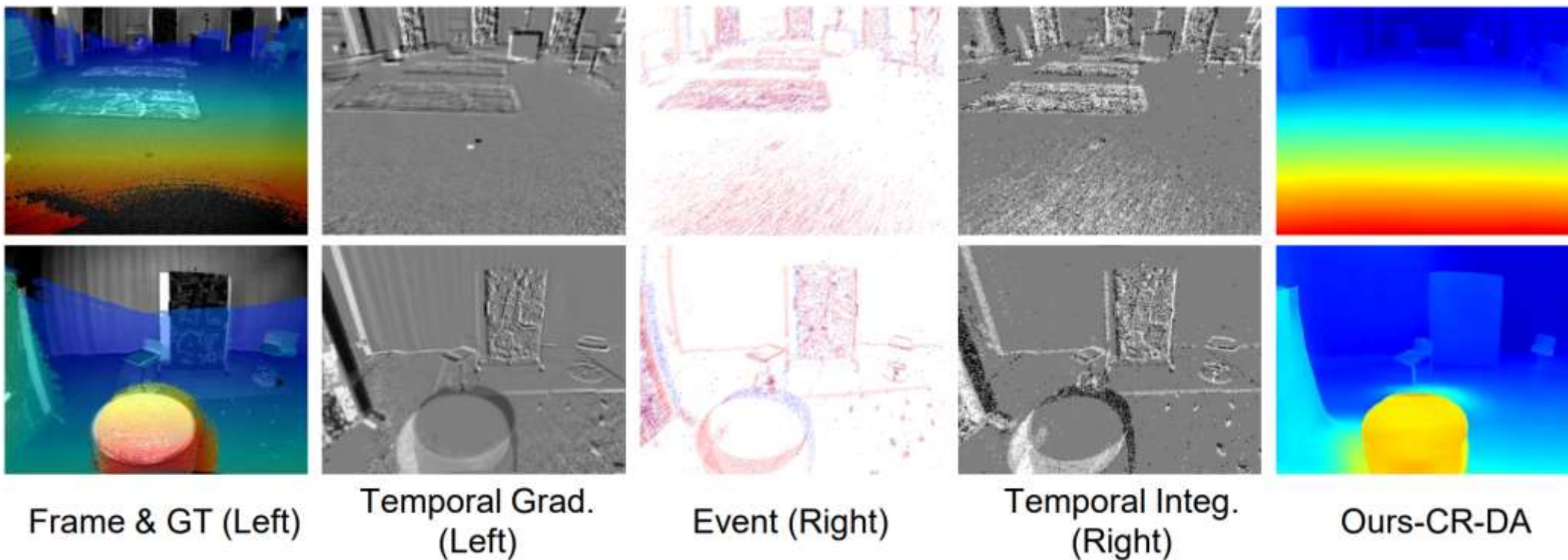




# Results



- Our method performs robustly in diverse scenarios and datasets.



# Thank You!

Lab page



<https://camera.pku.edu.cn>

Code available



<https://github.com/HYLZ-2019/ZEST>