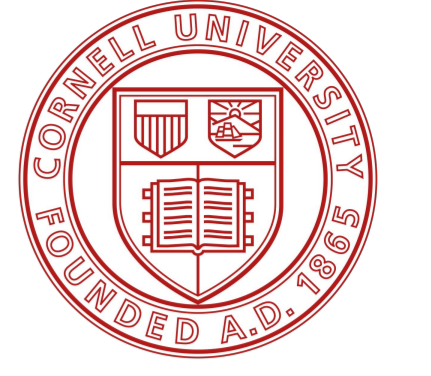


# The Caltech Fish Counting Dataset: A Benchmark for Multiple-Object Tracking and Counting

Justin Kay, Peter Kulits, Suzanne Stathatos, Siqi Deng, Erik Young, Sara Beery, Grant Van Horn, and Pietro Perona



## Counting Salmon in Sonar Video

Important application in conservation ecology: **how many salmon migrate upstream each year to spawn?**

Sonar video cameras are deployed in rivers as a **non-invasive and accurate** way to monitor salmon migration.

**Counting is currently performed manually** by technicians who watch video.

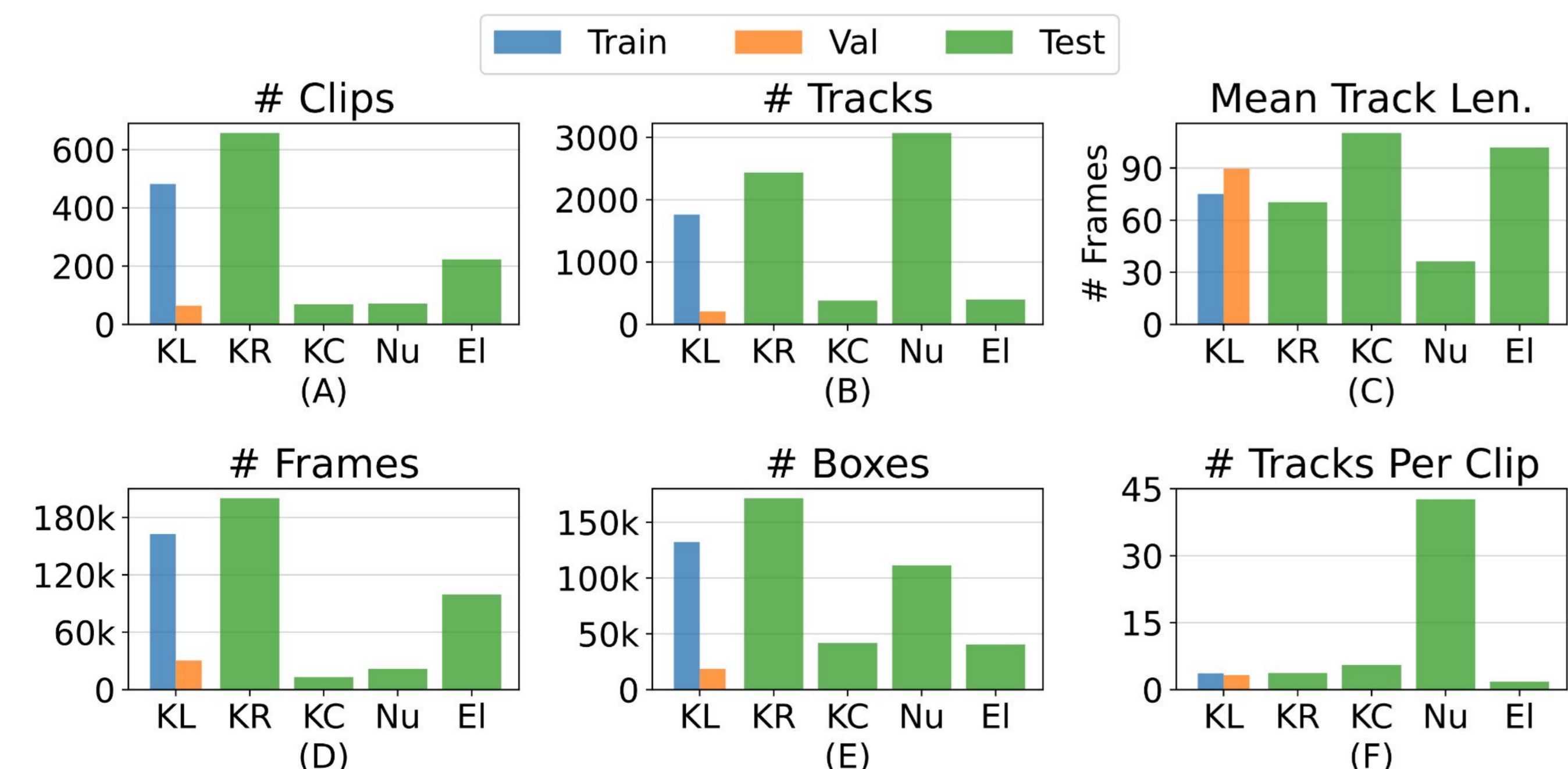
→ Automating counting is a **high-impact challenge** for computer vision.



## The Caltech Fish Counting Dataset

A large-scale dataset for **video object detection, multiple-object tracking, and video-based counting.**

- 1,567 video sequences (16.7 hours of video)
- 527,000 image frames
- 516,000 bounding boxes
- 8,254 object instances
- Test data from four out-of-distribution locations → study domain shift



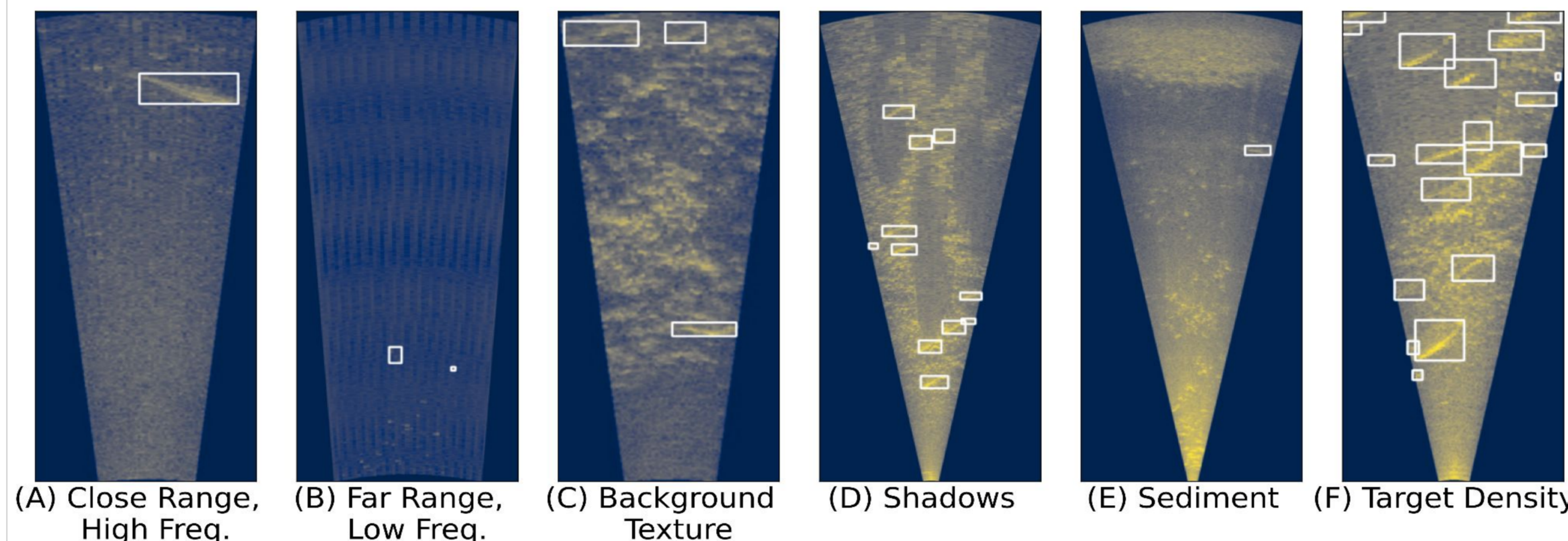
Five source locations: **KL**: Kenai River Left Bank; **KR**: Kenai River Right Bank; **KC**: Kenai River Minor Channel; **Nu**: Nushagak River; **EI**: Elwha River

## Data Challenges

### Domain shift

Each camera deployment presents different challenges to detection, tracking, and counting due to **location-specific environmental conditions.**

### Example frames and common challenges



### Performance degrades at OOD locations

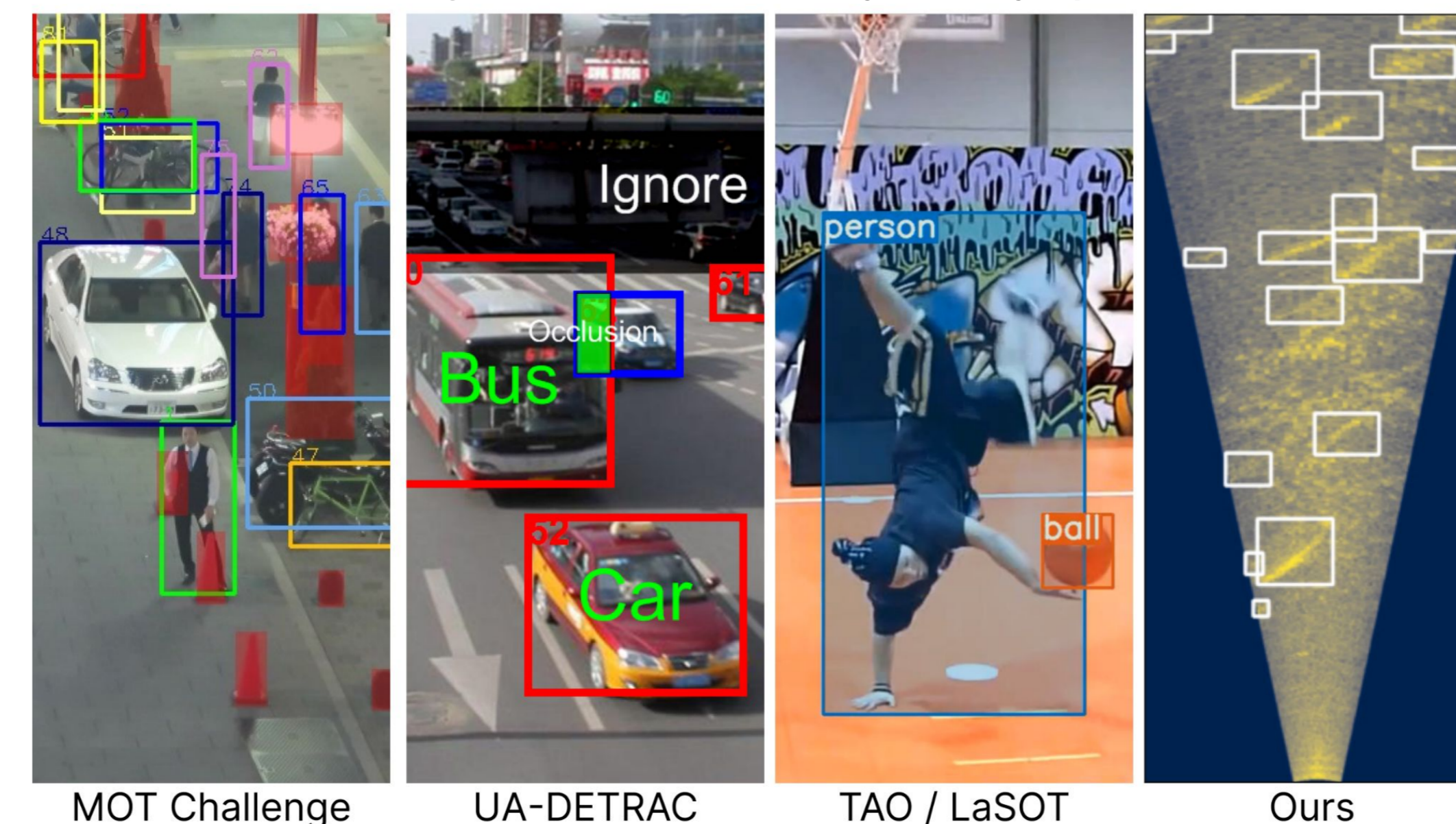
Loc	Baseline				
	AP <sub>@IoU=0.5</sub>	MOTA	IDF1	HOTA	nMAE
KL	66.4	44.9	66.7	49.2	4.9%
KR	57.7	-28.5	45.4	33.5	11.8%
KC	32.0	-60.8	35.6	30.9	53.0%
NU	70.6	30.2	60.8	44.4	14.0%
EL	39.9	-376.7	18.8	21.3	32.3%

### Low signal-to-noise data

- **Individuals look similar to background:** Detection not always possible in a single frame (need to incorporate temporal information).
- **Individuals look similar to each other:** Visual features ineffective for target re-identification.
- **Artifacts from sonar:** speckle noise, acoustic shadows, deterioration of signal at long-range, “ghost fish” (echoes that reflect off water surface)

### Comparison with other tracking datasets

In Caltech Fish Counting, trackers cannot rely heavily upon visual association



## Baseline Results

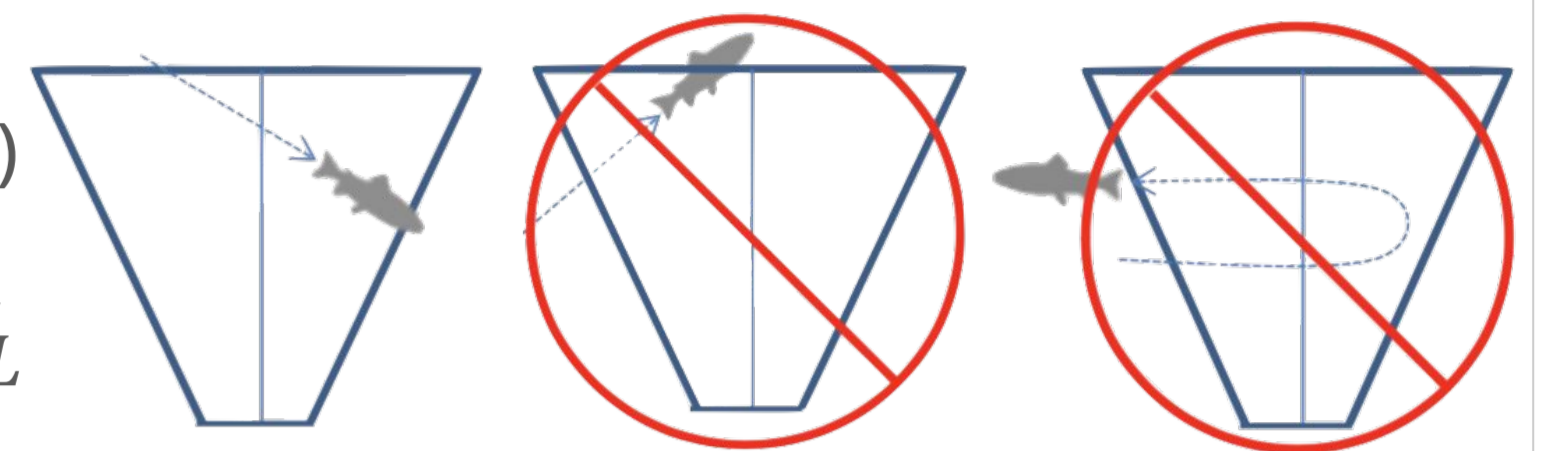
### Evaluation

- Detection: AP@IoU=0.5
- Tracking: MOTA, IDF1, HOTA
- Counting: normalized MAE (nMAE)

$$nMAE = \frac{\text{Abs. count err. at location } L}{\text{Ground truth count at location } L}$$

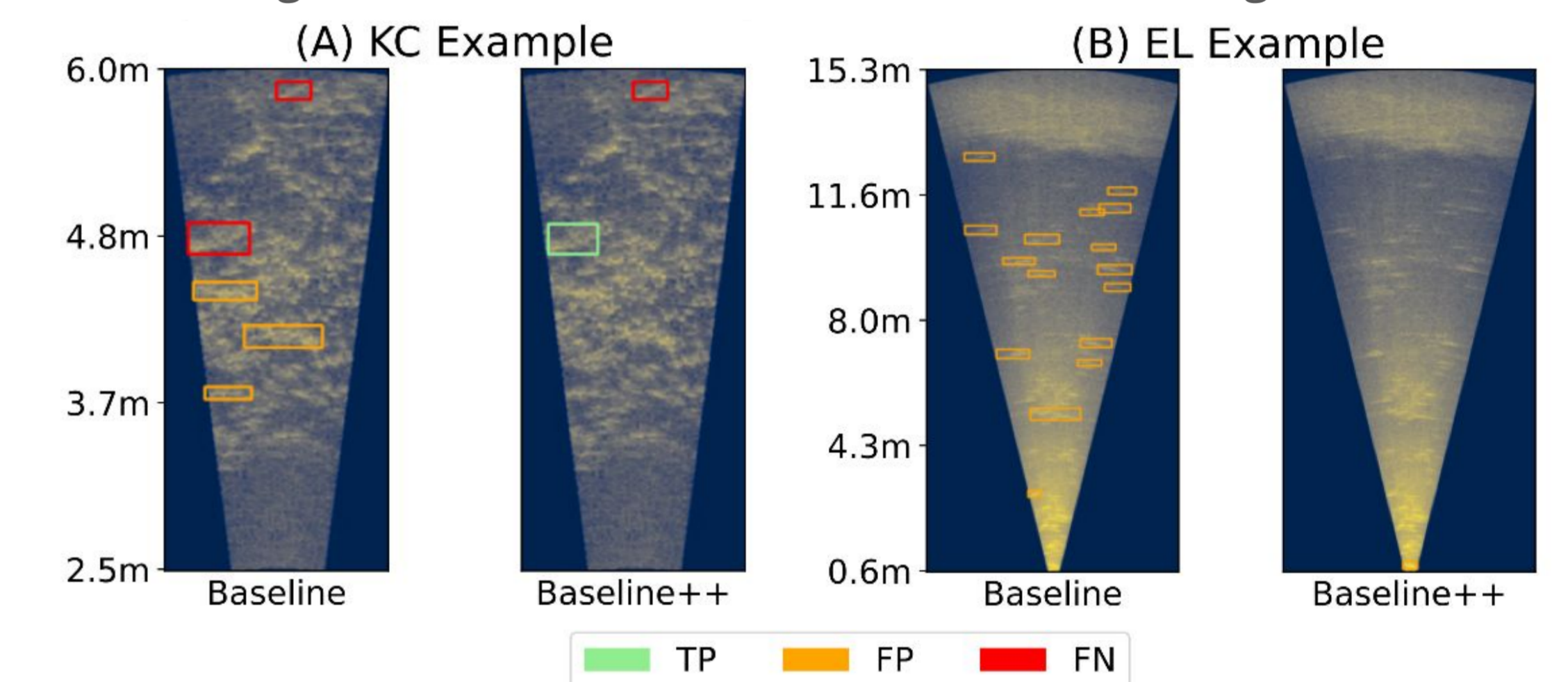
→ Target: < 10% nMAE

Crossing-line based counting  
Matches protocol used in the target application

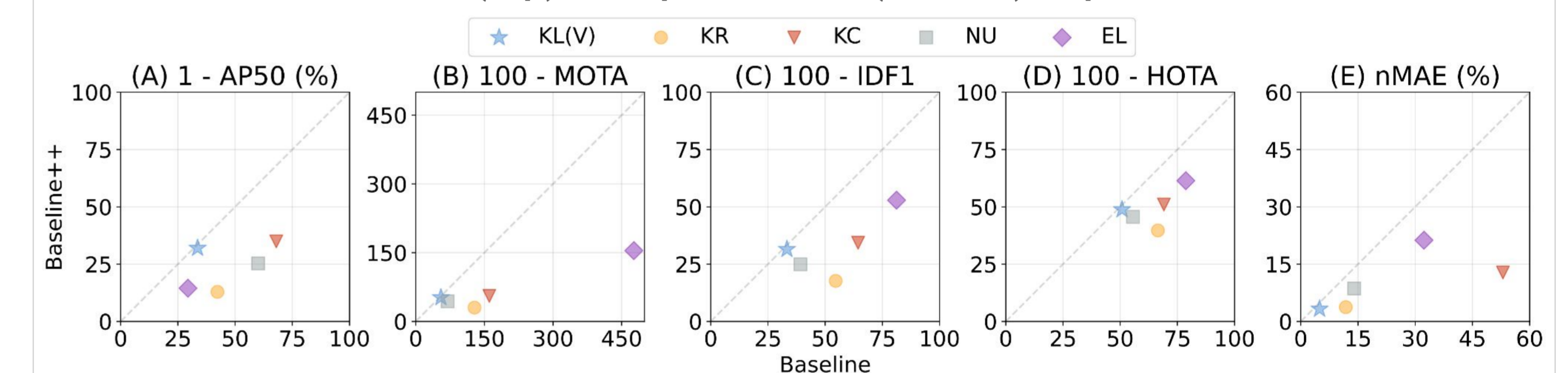


### Baseline and Baseline++

- YOLOv5m + SORT (Kalman filter + IoU-based association)
- “Baseline++” **incorporates temporal information** by augmenting input format w/ background subtraction and frame differencing

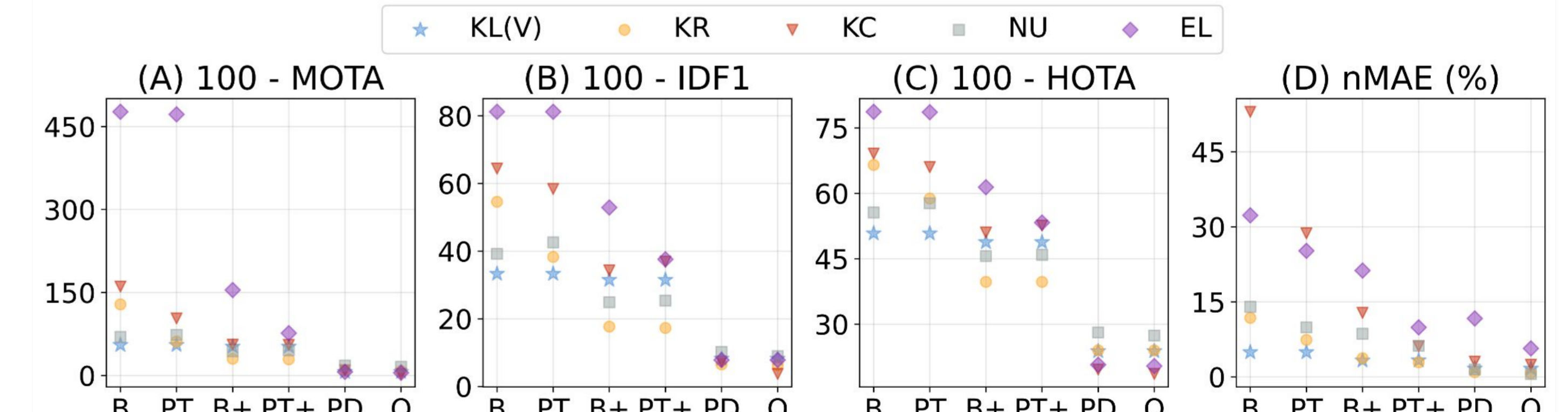


### Qualitative (top) and quantitative (bottom) improvements



### Upper bounds analysis

Compare **Baseline**, **Perfect Tracker**, **Baseline++**, **Perfect Tracker++**, **Perfect Detector**, **Oracle**



- Still **significant room for improvement** (13% nMAE at KC, 21% nMAE at EL)
- **Most effective path forward is detector improvements**

## Dataset Release



visipedia/caltech-fish-counting



Video example from dataset