# The 1st Visual Object Tracking Segmentation VOTS2023 Challenge Results

Matej Kristan, Jiri Matas, Martin Danelljan, Michael Felsberg, Hyung Jin Chang, Luka Čehovin Zajc, Alan Lukežič, Ondrej Drbohlav, Zhongqun Zhang, Khanh-Tung Tran, Xuan-Son Vu, Johanna Bjorklund, Christoph Mayer, Yushan Zhang, Lei Ke, Jie Zhao, Gustavo Fernandez

# Intro/Motivation

- The VOT formed in 2013 to support general object trackers

- "*Track any image region (including unknown instances or their parts) given a single training example.*"

- Research questions:
  Representations, Self-supervision, Robust localization…

- To better study these, the VOT challenges were restricted to single-target tracking, separating short-term and long-term tracking

# Intro/Motivation

- Ten challenges organized to explore various research questions



- The field has matured to a point where advancements expected by relaxing the restrictions

Visual object tracking segmentation challenge

VOTS

# The VOTS2023 challenge scope

- General object Short/Long-term, Single/Multi-target segmentation trackers

- Initialize on all targets in the first frame and report position in the rest



- Determine the target absence and redetect when it reappears

- Drifting off the target to background or another object is considered failure

# Per-target performance measures

## 5 tracking scenarios emerge:



Successfully localized | Tracker drift | Incorrectly predicted as absent | Incorrectly predicted as present | Correctly predicted as absent

Target present — Target absent

- **IoU as a standard measure** of agreement between prediction and GT

- Require IoU value definition for sc5

$$IoU_{sc5}=1.0$$

# Primary performance measure

- Performance summarized by the classical success w.r.t IoU plot (i.e., tracking quality plot)

- Success plot calculated individually for each target in each sequence and then averaged

- Primary measure: *Tracking quality Q* (area-under-the-curve)

Tracking quality (success) plot

Q = Area-under-the-curve

IoU threshold

# Auxiliary performance measures

- **Accuracy/Robustness** (@IoU=0.0 when target present)



*"Why did the tracker fail while target visible?"*

- N$_{ot}$R$_{eported}$E$_{rror}$ (NRE): % frames incorrectly predicted target absent

- D$_{rift}$R$_{ate}$E$_{rror}$ (DRE): % frames tracker drifted while predicting target present

*"How well is target absence determined?"*

- A$_{bsence}$D$_{etection}$Q$_{uality}$ (ADQ): % frames target correctly predicted absent

# VOTS2023 dataset

- Source: LaGOT[1], UTB180[2], TOTB[3], VOT-LT2021,VOT-LT2022, VOT-ST2022

- Selection criteria:

  - Sequences challenging for modern architectures

  - Properties: (i) visually-similar objects, (ii) substantial appearance changes, (iii) cluttered background, (iv) entering-exiting field-of-view

  - Diverse object and scene types (Air, Ground, Underwater)

  - Opaque as well as transparent objects

- Annotation: Segmentation masks

  - Include parts of objects as targets

[1] Mayer et al. ArXiv 2023; [2] Alawode et al. ACCV2022; [3] Fan et al. ICCV2021

# VOTS2023 dataset

- Stats: 144 sequences ; 341 targets ; 168 targets leave the FOV at least once

- Sequence properties:

  - min/max = 63/10.7k frames

  - On average 2.37 targets per sequence annotated

  - Median target absence: 18 frames

- To prevent overfitting:

  - Sequences + initialization frames GT publicly available.

  - GT of test frames sequestered, evaluation carried out on a dedicated server.

# VOTS2023 challenge results: 47 trackers tested

- Top trackers: DMAOT, HQTrack, MVOSTracker, $Dynamic_{DEAOT}$, seqtrack, DMNet, aot, MCMOT, rts_rts50_002, VAPT

- Dominant design choices:
  - Transformer-based
  - Single-stage ST1/LT0 trackers
  - Same architecture used for frame-to-frame localization and re-detection

# VOTS2023 challenge quality of submissions

- Baseline 1: Independent STARKs[1]  (47% in Q w.r.t. top tracker)

  80% of submissions outperform it

- Baseline 2: VOT2022 winner AOT[2]

  13% (top 6 trackers) outperfrom it



[1] Yan, et al. ICCV2021; [2] Yang, et al. NeurIPS 2021

# VOTS2023 challenge results

- Top performer DMAOT: Extends the VOT2022 winner AOT 🏆
  - *Swin transformer backbone ; Separates long-term and short-term target templates; gated propagation module for visual embeddings; NCV motion model*

- Very good Acc=0.751 & Rob=0.795 (localizes the target 80% of the time)

- Very low drifting (DRE=7%),

- Low false absence prediction (NRE=14%)

- Good target absence prediction: in ADQ=73% cases



DMAOT

# VOTS2023 challenge results

- The top-performer in Q (DMAOT) strikes a good balance in Acc/Rob

- Top robustness: DMNet (Rob=0.86) vs (DMAOT Rob = 0.795)

  - Reason might be the use of optimal transport formulation in segmentation/localization

- Top accuracy: SeqTrack

  - Bounding box tracker with SAM[1] segmentation

  - Care taken when to accept the SAM[1] result

[1]Kirillov, et al., Segment Anything, 2023

# VOTS2023 challenge results

- Q @ low thresholds indicates robustness

- Q @ mid-to-high thresholds indicates accuracy

Clusters:

- @ low thresholds

  - Transformer feature extraction backbones

- @ medium-to-high thresholds

  - Careful use of SAM[1] for segmenting targets (mask or box refinement)

[1]Kirillov, et al., Segment Anything, 2023

VOTS2023 challenge Winners:

DMAOT by: Yangming Cheng, Zongxin Yang, Yuanyou Xu, Xiaodi Li, Jiahao Li, Yi Yang, Yueting Zhuang

"Decoupled Memory AOT"

VOTS2023 challenge Spotlight:

DMNet by: Yinchao Ma, Wangkai Li, Dawei Yang, Rui Sun, Qianjin Yu, Fei Wang, Tianzhu Zhang

"Dynamic Matching Network"

Winners & Spotlight talks in Session II @10:45

# Summary

- New challenge: General Short/long-term, Single/Multi-target segmentation

- New performance measures, dataset, toolkit and eval server

- Evaluation server open for post-challenge evaluation





VOTS benchmark

- Similarly to VOT2022, evidence indicates remarkable robustness of segmentation trackers vs bbox trackers

- Encourage evaluation of bounding box trackers on the VOTS2023 benchmark (Robustness)

# Thanks

- ## The VOTS2023 committee



M. Kristan     J. Matas     M. Danelljan     M. Felsberg     H. J. Chang     L. Čehovin Z.     A. Lukežič     O. Drbohlav     Z. Zhang
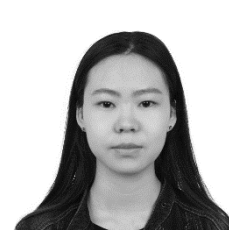
T. Tran     Xuan-Son Vu     Johanna Björklund     C. Mayer     Z. Yushan     Lei Ke     G. Fernandez

- ## Everyone who participated or contributed

Noor Al-Shakarji38, Dong An20, Michael Arens15, Stefan Becker15, Goutam Bhat3, Sebastian Bullinger15, Antoni B. Chan11, Shijie Chang13, Hanyuan Chen14, Xin Chen13, Yan Chen19, Zhenyu Chen13, Yangming Cheng42, Yutao Cui29, Chunyuan Deng16, Jiahua Dong32, Matteo Dunnhofer41, Wei Feng34, Jianlong Fu27, Jie Gao19, Ruize Han34, Zeqi Hao13, Jun-Yan He14, Keji He20, Zhenyu He18, Xiantao Hu17, Kaer Huang25, Yuqing Huang18, Yi Jiang9, Ben Kang13, Jin-Peng Lan14, Hyungjun Lee30, Chenyang Li14, Jiahao Li42, Ning Li17, Wangkai Li39, Xiaodi Li42, Xin Li31, Pengyu Liu13, Yue Liu23, Huchuan Lu13, Bin Luo14, Ping Luo33, Yinchao Ma39, Deshui Miao18, Christian Micheloni41, Kannappan Palaniappan38, Hancheol Park30, Matthieu Paul3, HouWen Peng26, Zekun Qian34, Gani Rahmon38, Norbert Scherer-Negenborn15, Pengcheng Shao23, Wooksu Shin30, Elham Soltani Kazemi38, Tianhui Song29, Rainer Stiefelhagen24, Rui Sun39, Chuanming Tang37, Zhangyong Tang23, Imad Eddine Toubal38, Jack Valmadre35, Joost van de Weijer12, Luc Van Gool3, Jash Vira35, St`ephane Vujasinovi´c15, Cheng Wan16, Jia Wan8, Dong Wang13, Fei Wang39, Feifan Wang34, He Wang23, Limin Wang29, Song Wang40, Yaowei Wang31, Zhepeng Wang25, Gangshan Wu29, Jiannan Wu33, Qiangqiang Wu11, Xiaojun Wu23, Anqi Xiao20, Jinxia Xie17, Chenlong Xu17, Min Xu10, Tianyang Xu23, Yuanyou Xu42, Bin Yan13, Dawei Yang39, Ming-Hsuan Yang36, Tianyu Yang22, Yi Yang42, Zongxin Yang42, Xuanwu Yin28, Fisher Yu3, Hongyuan Yu28, Qianjin Yu39, Weichen Yu10, YongSheng Yuan13, Zehuan Yuan9, Jianlin Zhang37, Lu Zhang13, Tianzhu Zhang39, Guodongfang Zhao21, Shaochuan Zhao23, Yaozong Zheng17,19, Bineng Zhong17, Jiawen Zhu13, Xuefeng Zhu23, Yueting Zhuang42, ChengAo Zong13, and Kunlong Zuo28

- ## VOTS2023 sponsors:



University of Ljubljana
Faculty of Computer and
Information Science

eyedea

UNIVERSITY OF
BIRMINGHAM