

HOW VIDEO OBJECT TRACKING IS AFFECTED BY IN-CAPTURE DISTORTIONS?

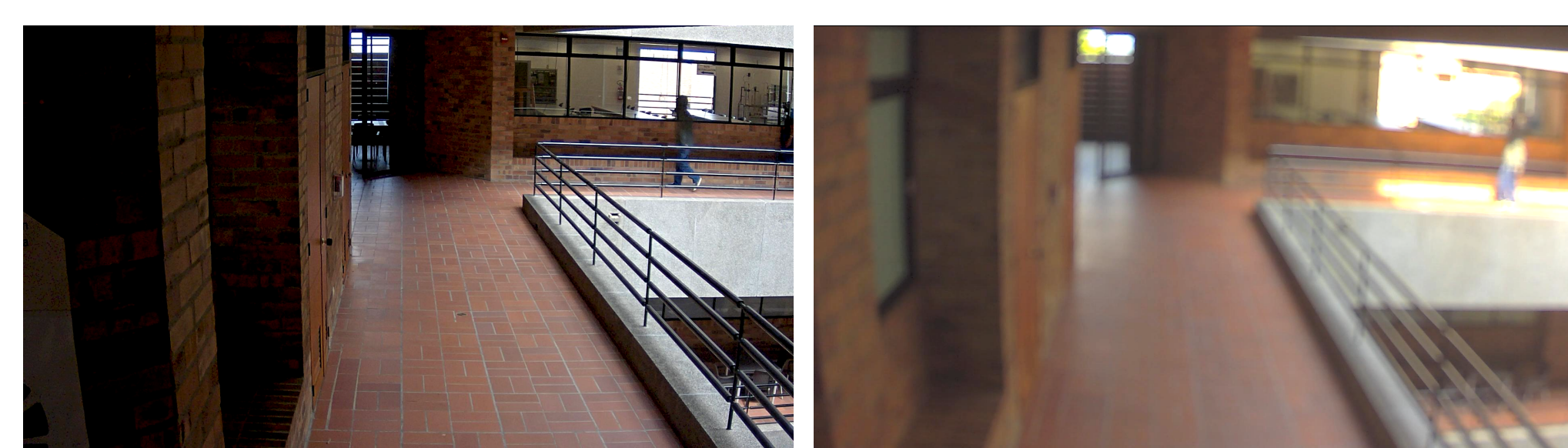
Roger Gomez Nieto¹, Hernan Dario Benitez Restrepo¹, Ivan Cabezas²

¹Pontificia Universidad Javeriana - Cali, Colombia. {roger.gomez, hbenitez}@javerianacali.edu.co

²Universidad de San Buenaventura- Cali, Colombia. imcabezas@usbcali.edu.co

SUMMARY

- Video Object Tracking -VOT- in realistic scenarios is a difficult task. Image factors such as occlusion, clutter, confusion, object shape, and zooming have an impact on video tracker methods performance.
- There is a lack of a detailed study analyzing performance on videos with authentic in-capture distortions. Such a study requires a database with videos containing distortions in a controlled and quantifiable way.



(a) Pristine Indoor (b) Distorted Indoor



(c) Pristine Outdoor (d) Distorted Outdoor

Figure: Examples of pristine and distorted images within indoor and outdoor environments.

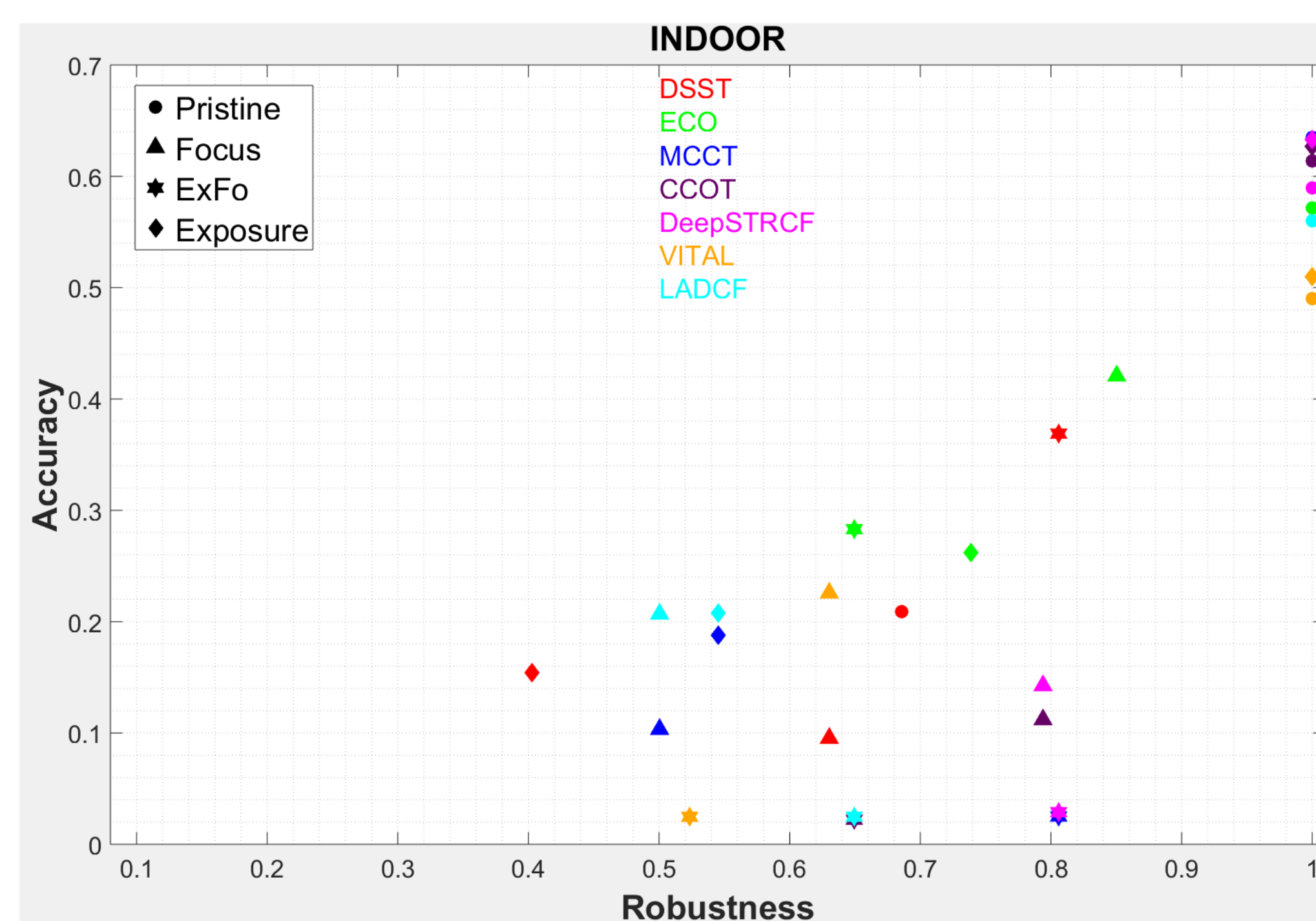


Figure: A-R plot for VOT in an indoor environment with pristine and distorted videos with the same activity.

MATERIALS AND METHODS

We used the A-R plot performance measure [1] for analyzing tracker accuracy and robustness. DSVD contains 537 videos affected by in-capture distortions, acquired by four different surveillance cameras. DSVD contains real-world surveillance scenes such as people walking alone, meeting, fighting, passing out, leaving a package in a public place, prowling, and being robbed. The videos have an equal rate I/P frames: 10 fps. The frame size is FHD (1920× 1080), the color space is three RGB channels and the exposure variation range is $\{\frac{1}{480}, \frac{1}{120}\}$ seconds. The video dataset also contains H.264/AVC compression post-capture distortions at three different bitrates. The three different bitrates are 4700, 1800 and 1200 kbps.

Each tracker was executed 30 times on each sequence, considering stochastic processes. We tested the trackers: CCOT, MCCT, ECO, DSST, DeepSTRCF, LADCF, and VITAL [2].

RESULTS

The trackers were tested in 15 scenes containing in-capture distortions such as lack of exposure, out of focus and out-of-focus concurrently with lack of exposure at indoor and outdoor environments.

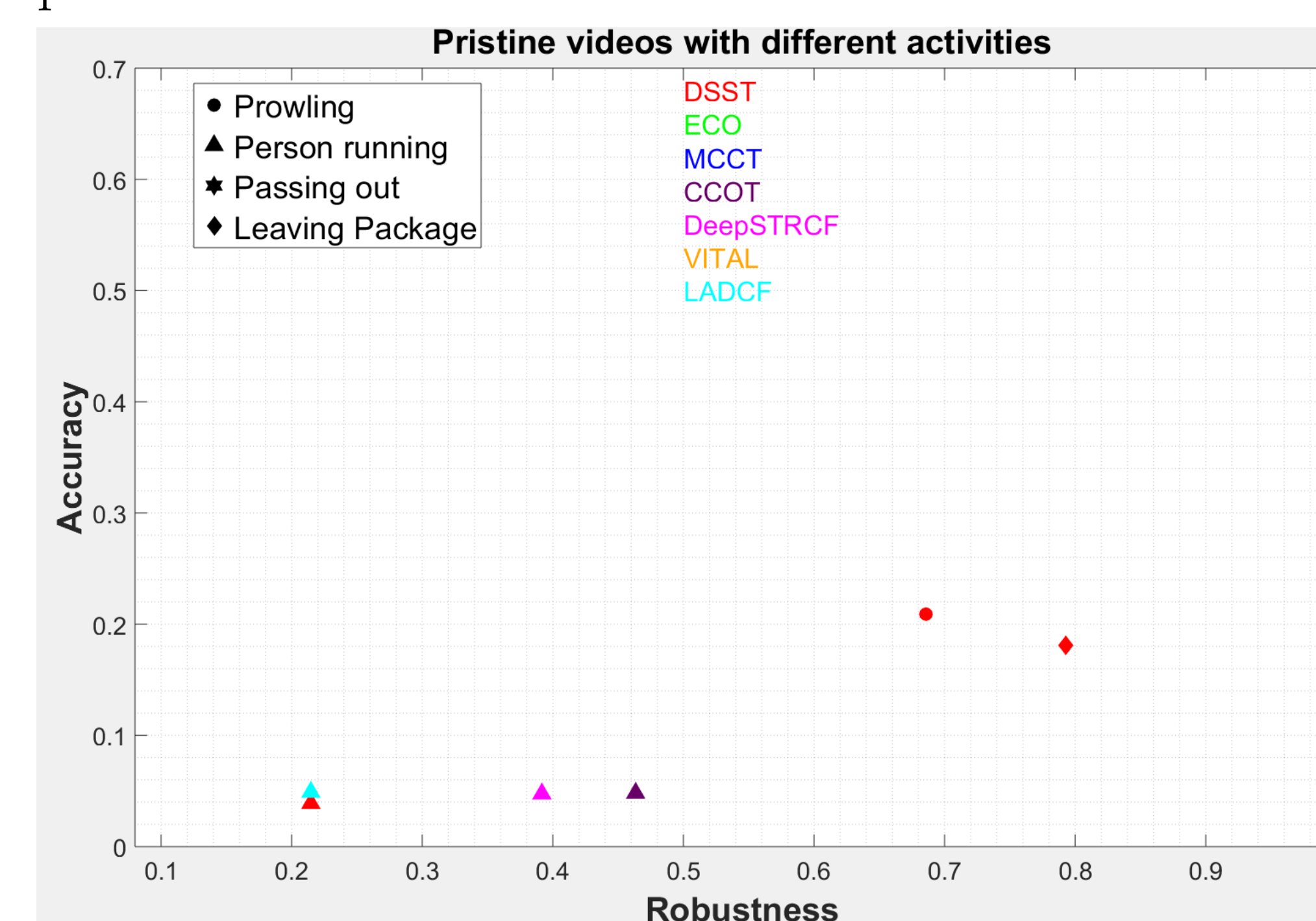


Figure: A-R plot for VOT within an indoor environment in pristine videos (same conditions for all) and different activity.

CONCLUSION

- DSVD with 537 surveillance videos containing different levels of authentic distortions such as low exposure and out-of-focus. DSVD can be seen as a solid starting point to study the influence of distortions on video tracker performance.
- In-capture distortions severely hamper VOT methods performance in a non intuitive way. It was shown by assessing seven state-of-the-art trackers using the A-R plot performance measure on DSVD.
- In practice, no specific type of distortion consistently generated the worst performance in all scenes, neither affected all trackers in the same way.

REFERENCES

- [1] L Cehovin, A Leonardis, and M Kristan. Visual Object Tracking Performance Measures Revisited. *Image Processing, IEEE Transactions on*, 25(3):1261–1274, 2016.
- [2] Matej Kristan and Ales Leonardis et al. The sixth visual object tracking vot2018 challenge results, 2018.

IMPORTANT RESULT

Distorted Surveillance Video Dataset (DSVD): **537 surveillance videos containing different levels of authentic distortions**. Assessment of seven state-of-the-art trackers on this dataset, demonstrating that in-capture distortions severely hamper VOT methods performance **in a non intuitive way**.

ACKNOWLEDGEMENTS

The authors acknowledge the funding provided by COLCIENCIAS and Pontificia Universidad Javeriana. The authors would like to thank NVIDIA Corporation for the donation of a TITAN XP GPU used in these experiments. The authors would also like to acknowledge the grant provided by *Comision Fulbright Colombia*

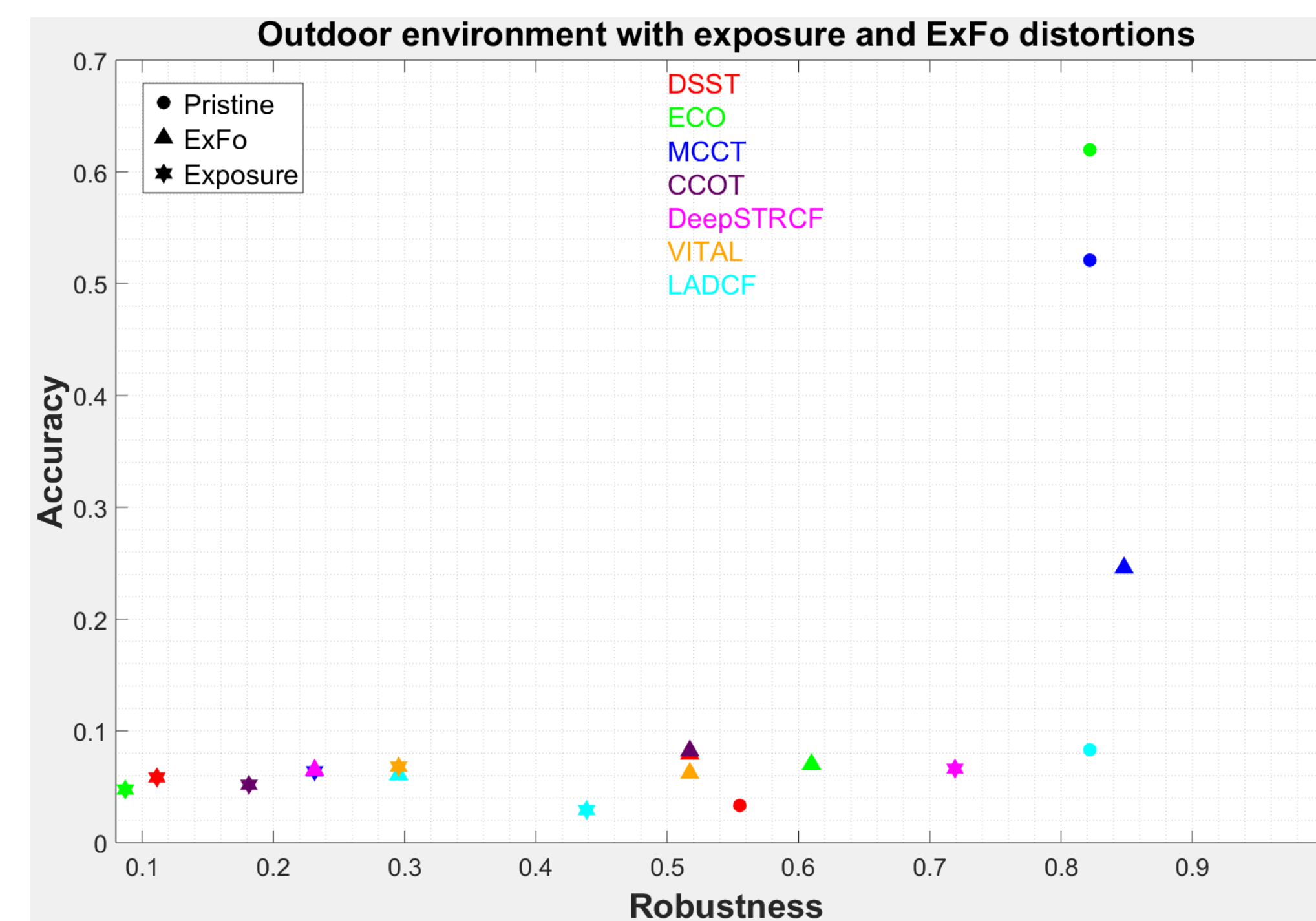


Figure: A-R plot for VOT in an outdoor environment with pristine and distorted videos with the same activity.

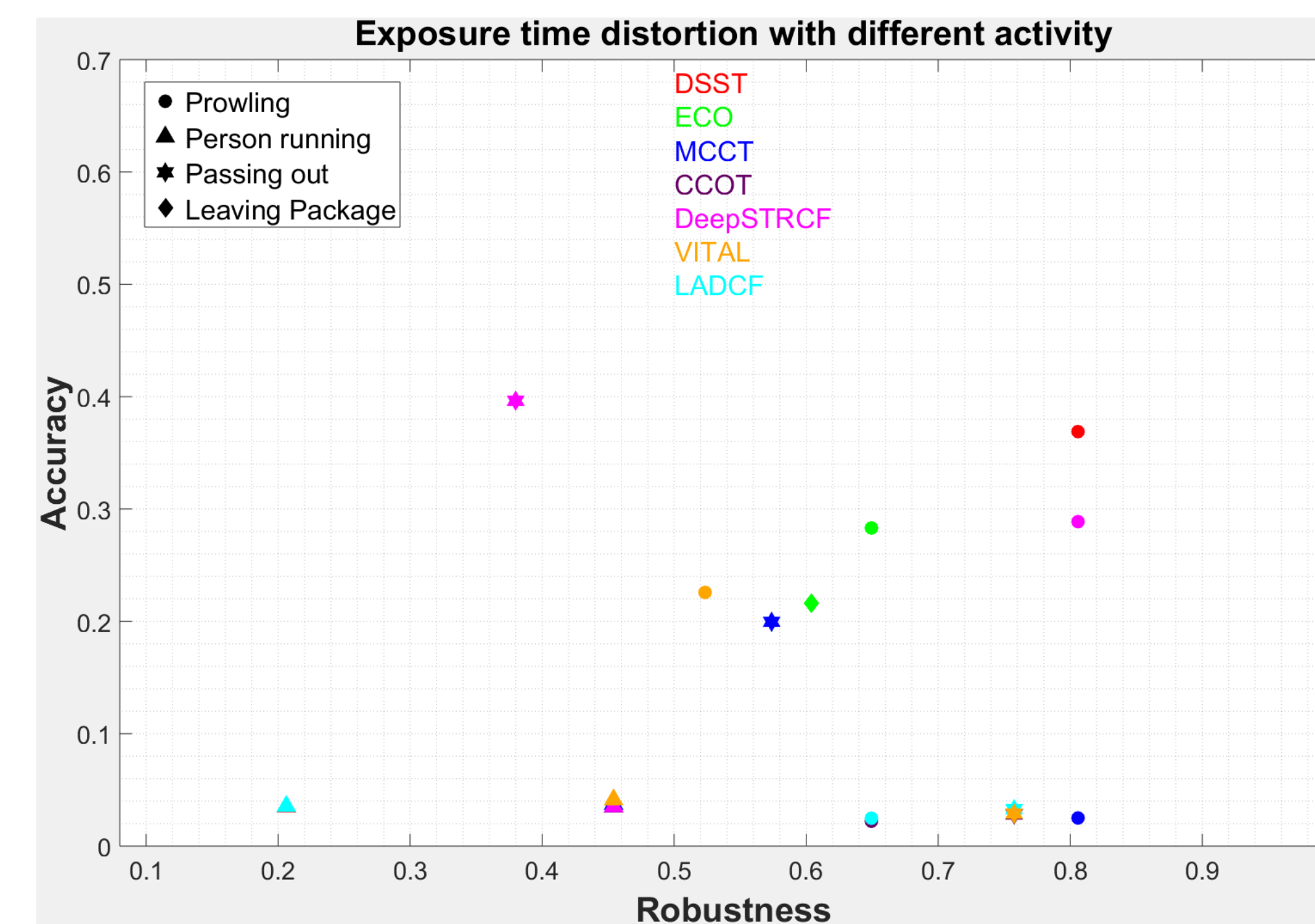


Figure: A-R plot for VOT with exposure time distortion in the same level and different activity in video.