



Deep Learning of Invariant Features via Fixations in Video

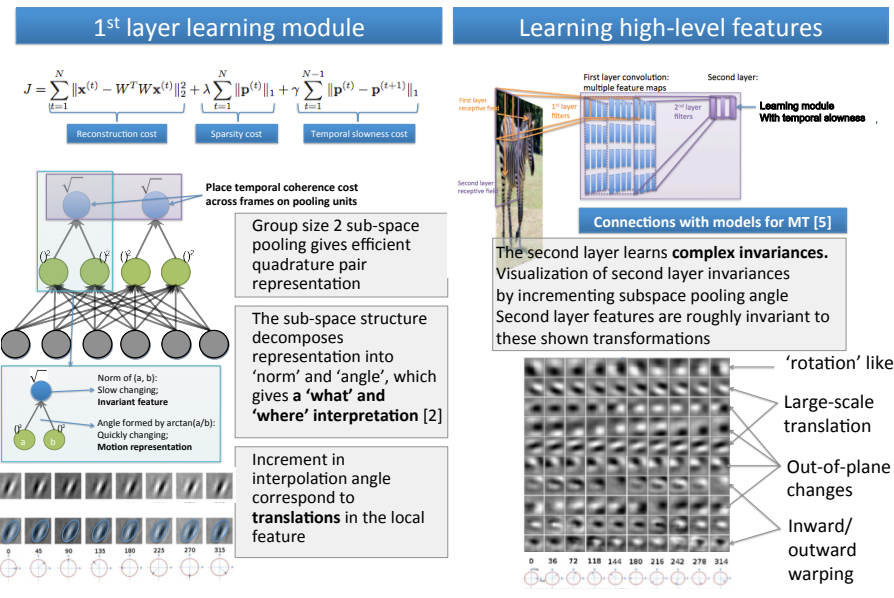
Will Y. Zou, Shenghuo Zhu, Andrew Y. Ng and Kai Yu

NEC Laboratories
America
Relentless passion for innovation

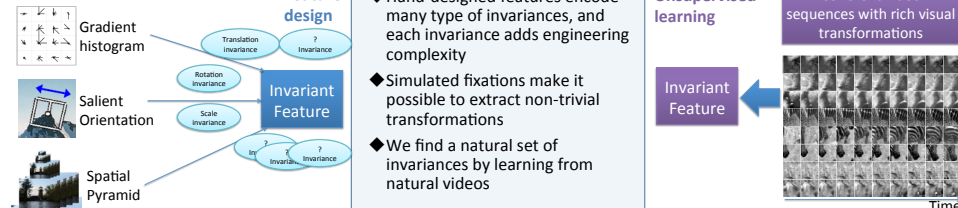
1 Abstract

We apply salient feature detection and tracking in videos to simulate fixations and smooth pursuit in human vision. With tracked sequences as input, a hierarchical network of modules learns invariant features using a temporal slowness constraint. The network encodes invariance which are increasingly complex with hierarchy. Although learned from videos, our features are spatial instead of spatial-temporal, and well suited for extracting features from still images. We applied our features to four datasets (COIL-100, Caltech 101, STL-10, PubFig), and observe a consistent improvement of 4% to 5% in classification accuracy. With this approach, we achieve state-of-the-art recognition accuracy 61% on STL-10 dataset.

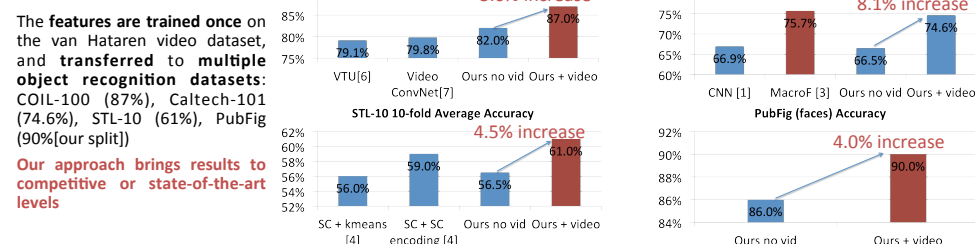
3 Learn hierarchically invariant features, unsupervised



2 Motivation



4 Results

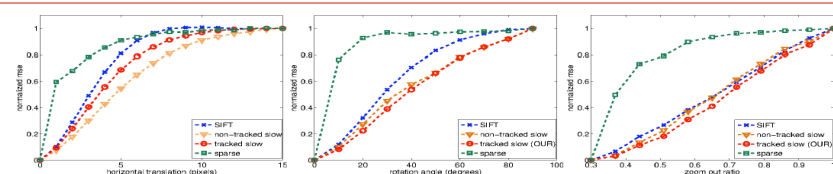


5 Invariance

Translation

Rotation

Zoom



6 References

- [1] K. Jarrett, K. Kavukcuoglu, M. Razento and Y. LeCun. What is the best Multi-Stage Arch. for Obj. Recog? ICCV 2009.
- [2] C. Cadieu, B. Olshausen. Learning Transformational Invariants from Natural Movies. NIPS 2010
- [3] Y. Boureau, F. Bach, Y. LeCun, and J. Ponce. Learning mid-level features for recognition. CVPR 2010
- [4] A. Coates, A. Y. Ng. The Importance of Encoding Versus Training with Sparse Coding and Vector Quantization. ICML 2011
- [5] N. Rust, V. Manté, E. P. Simoncelli, and J. A. Movshon. How mt cells analyze the motion of visual patterns. Nature Neuroscience, 2006
- [6] H. Wersing and E. Kröner. Learning optimized features for hierarchical models of invariant object recognition. Neural Comp. 2003
- [7] H. Mobahi, R. Collobert and J. Weston. Deep learning from temporal coherence in video. ICML 2009

<http://ai.stanford.edu/~wzou/>