

On Composition of Transformations for Self-Supervised Learning

Mandela Patrick*, Yuki M. Asano*, Polina Kuznetsova, Ruth Fong, João F. Henriques,
Geoffrey Zweig, Andrea Vedaldi

ICCV 2021

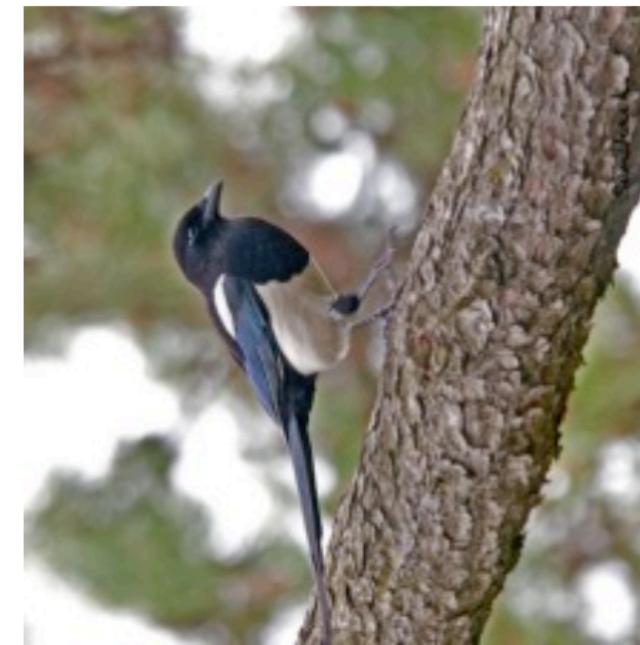
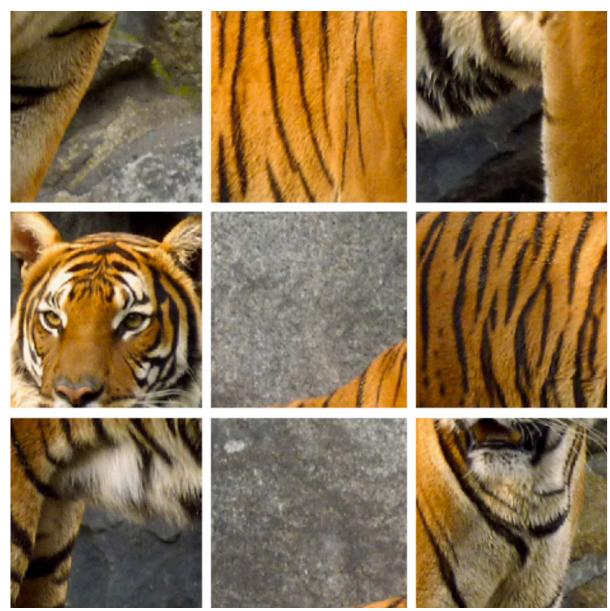
* equal contribution



UNIVERSITY OF
OXFORD

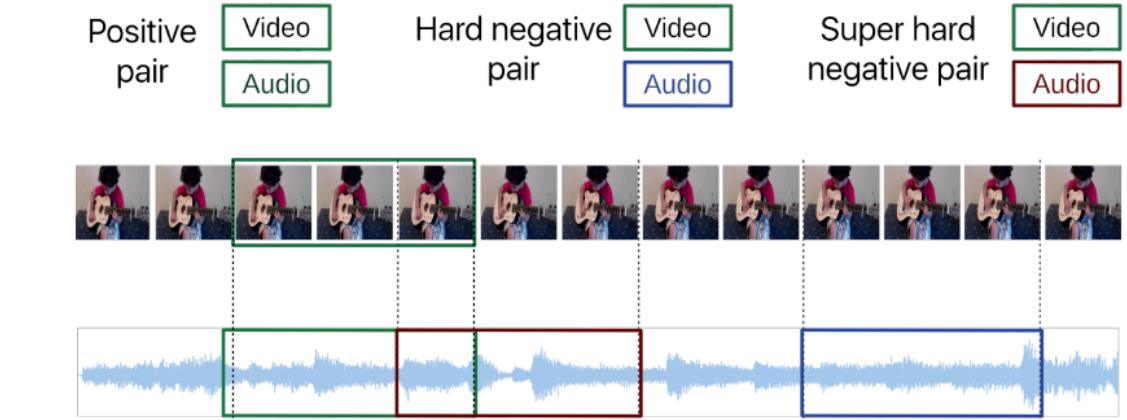
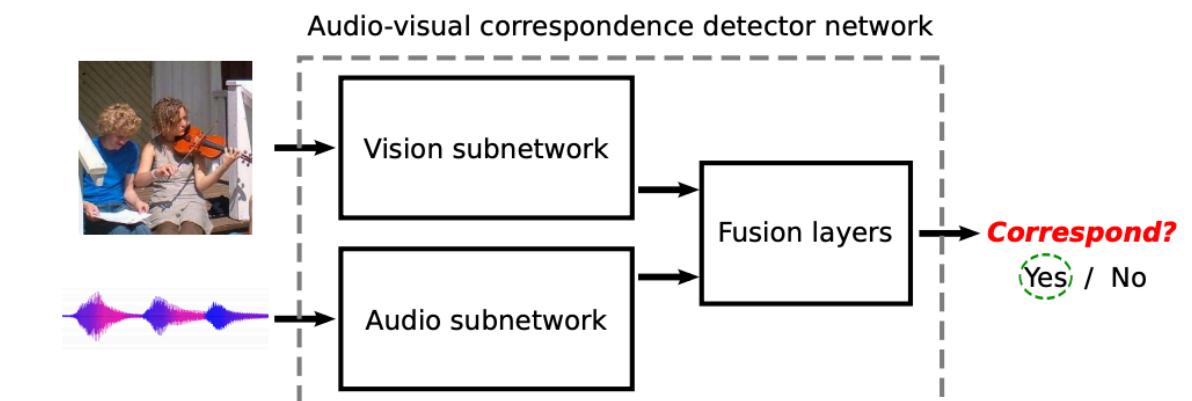
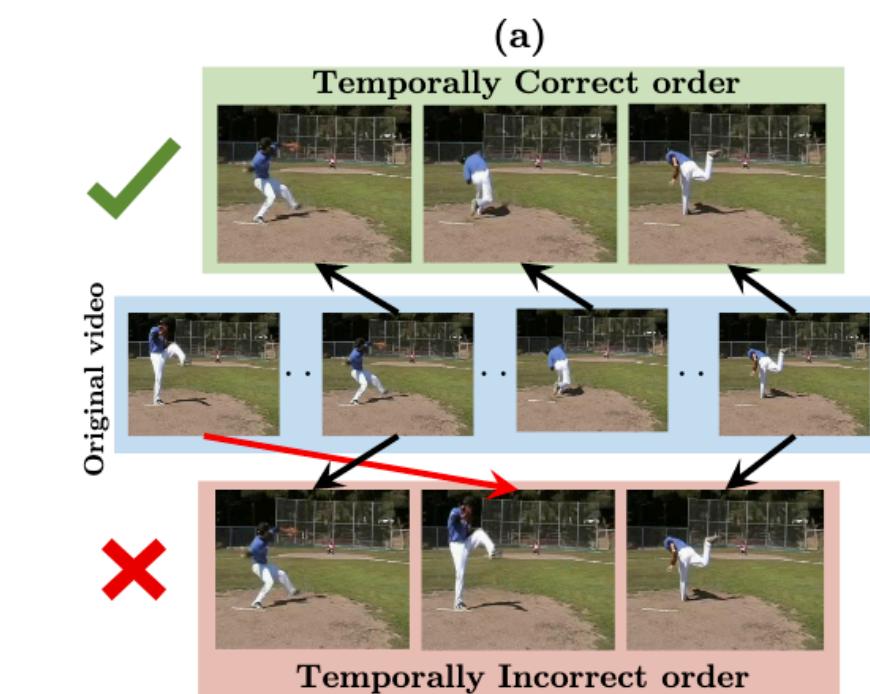
FACEBOOK

Self-supervision: a pretext zoo



Images:

Jigsaw, RotNet, Colorisation



Videos: two more dimensions
(time and modality)
Shuffle&Learn, L3, AVTS

Noise contrastive learning

$$f\left(\begin{array}{c} \text{Image of a cat with blue eyes} \end{array}\right) = f\left(\begin{array}{c} \text{Image of a cat with blue eyes} \end{array}\right) \neq f\left(\begin{array}{c} \text{Image of a white dog} \end{array}\right)$$

Key idea:

features should encode image's core information.

learn this by comparing augmentations against other images.

Examples: *NPID, MoCo, CMC, SimCLR..*

[Wu et al., CVPR 2018; He et al., CVPR 2020; Tian et al., ECCV 2020; Chen et al., ICML 2020]

Most pretext tasks

$$f\left(\begin{array}{c} \text{Image of a cat} \end{array} \right) = \text{Clockwise 90 degrees rotation}$$

e.g. RotNet

Pretext tasks = *Specifying* invariance and distinctiveness to *selected* transformations.

Invariance vs distinctiveness

Should the representation discount or distinguish time reversal?



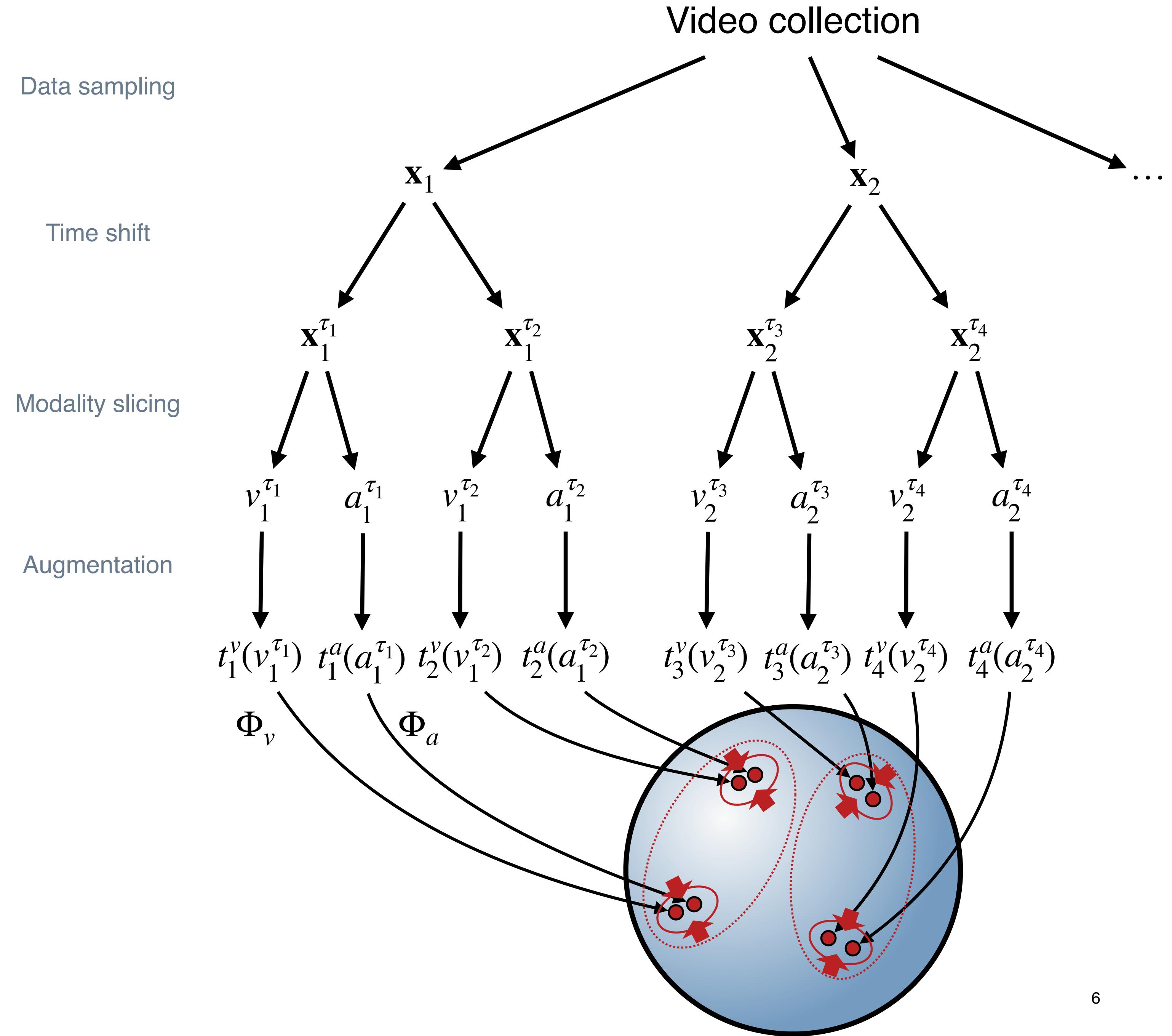
?
=



GDT: Generalised data transformations

GDTs express all aspects of these methods as transformations.

Allow to specify either invariance or distinctiveness



GDT Transformations

In our work, we explore the following transformations as learning hypotheses:

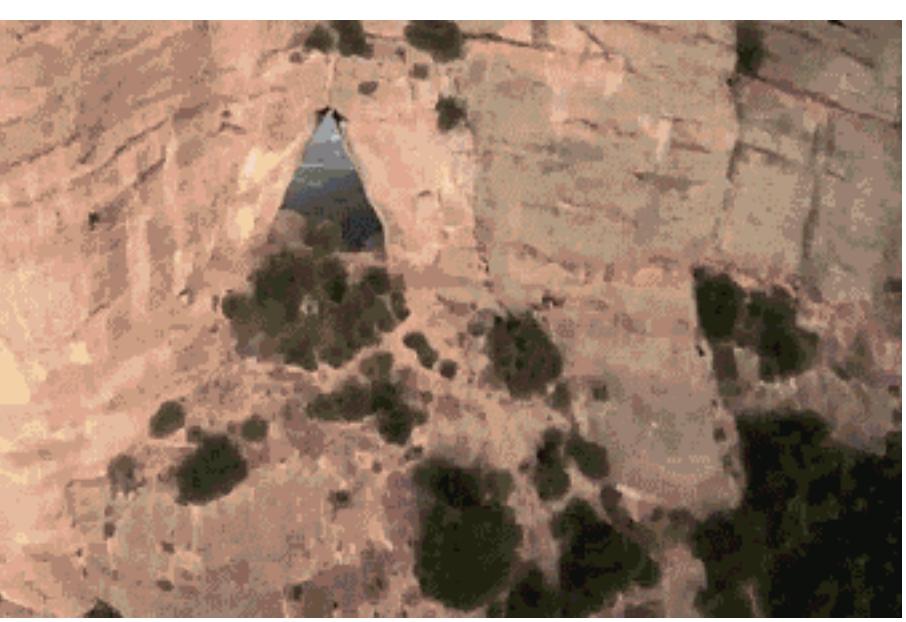
1. Sample Distinctiveness



2. Time Reversal

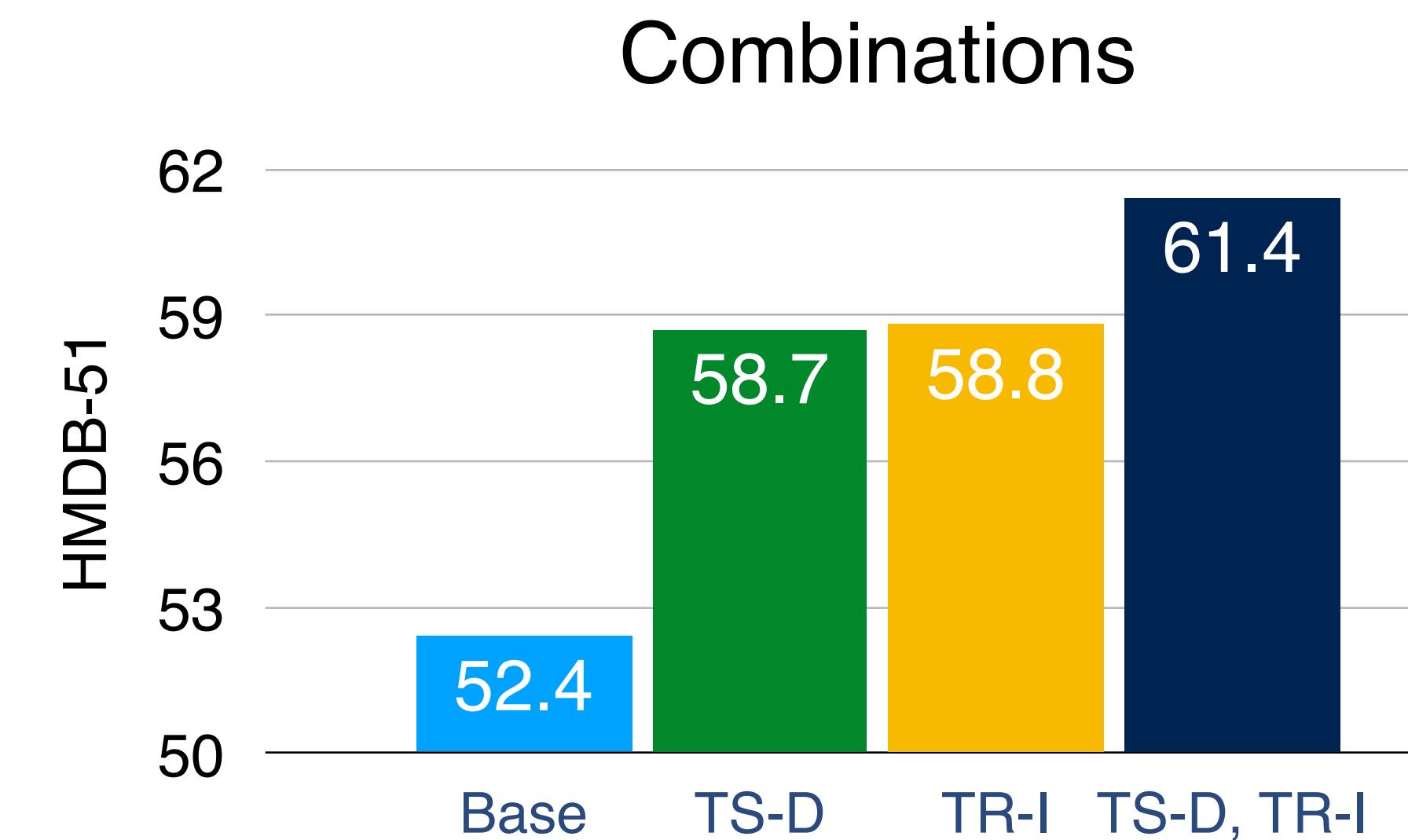
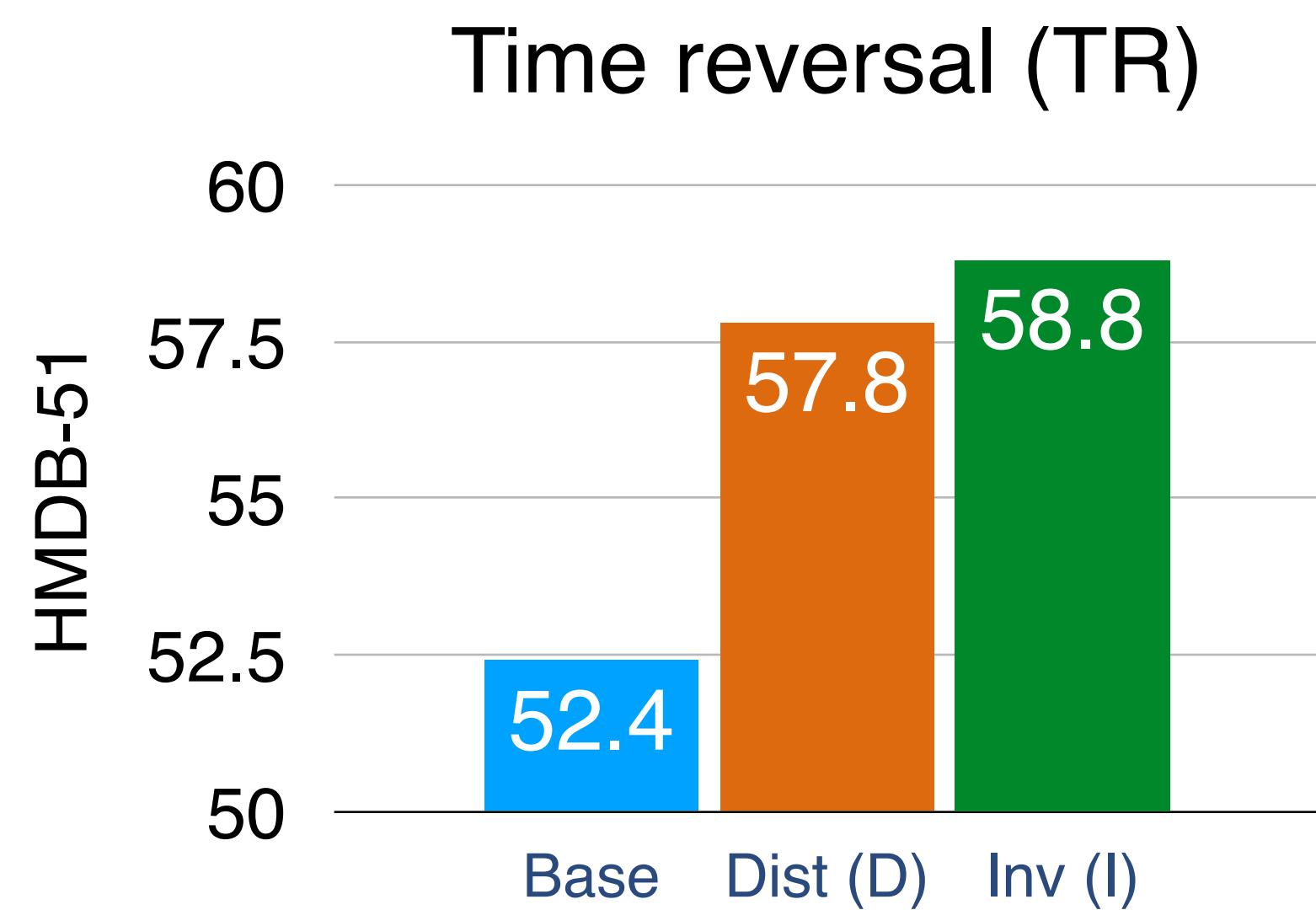
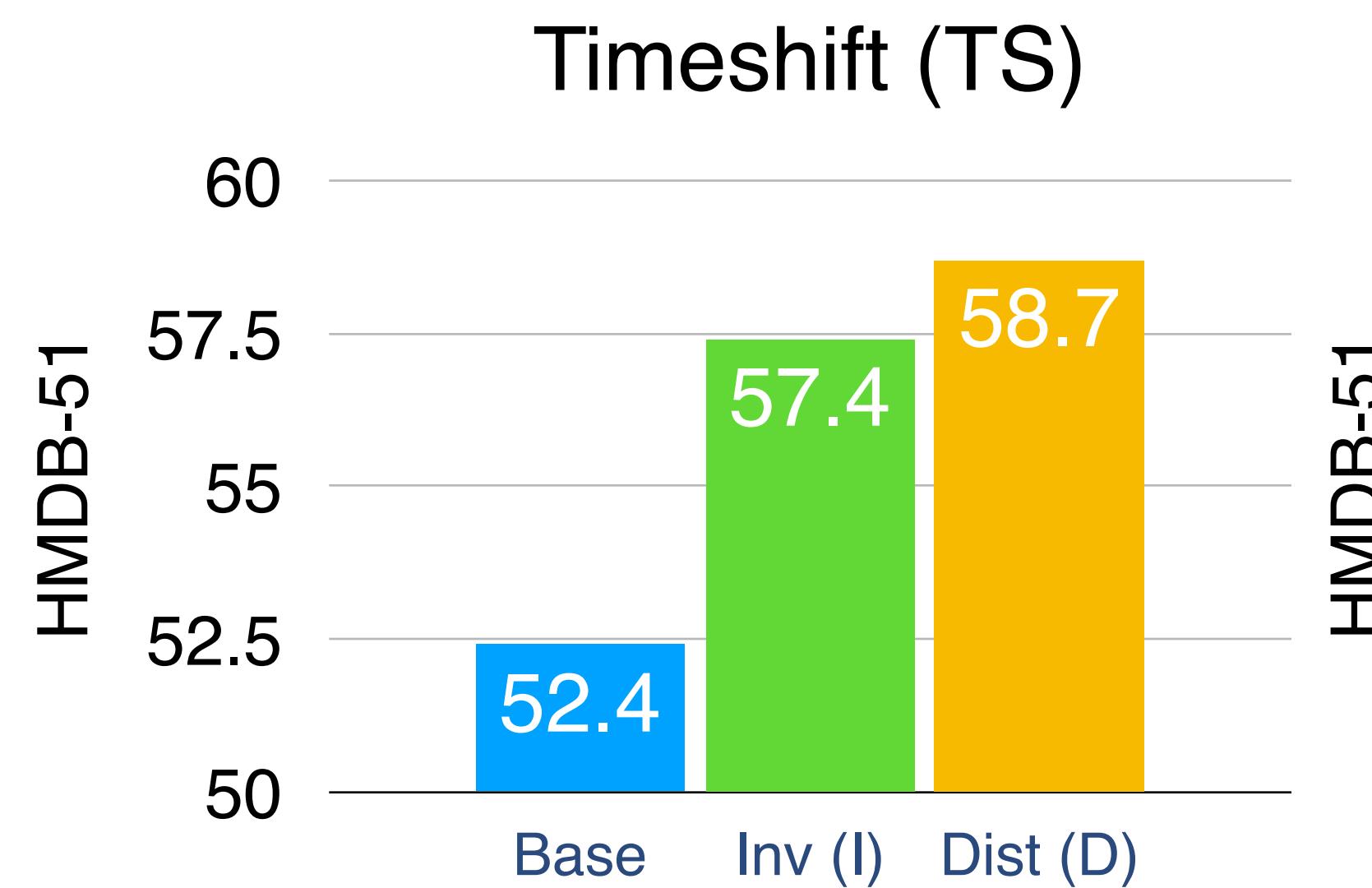


3. Time Shift

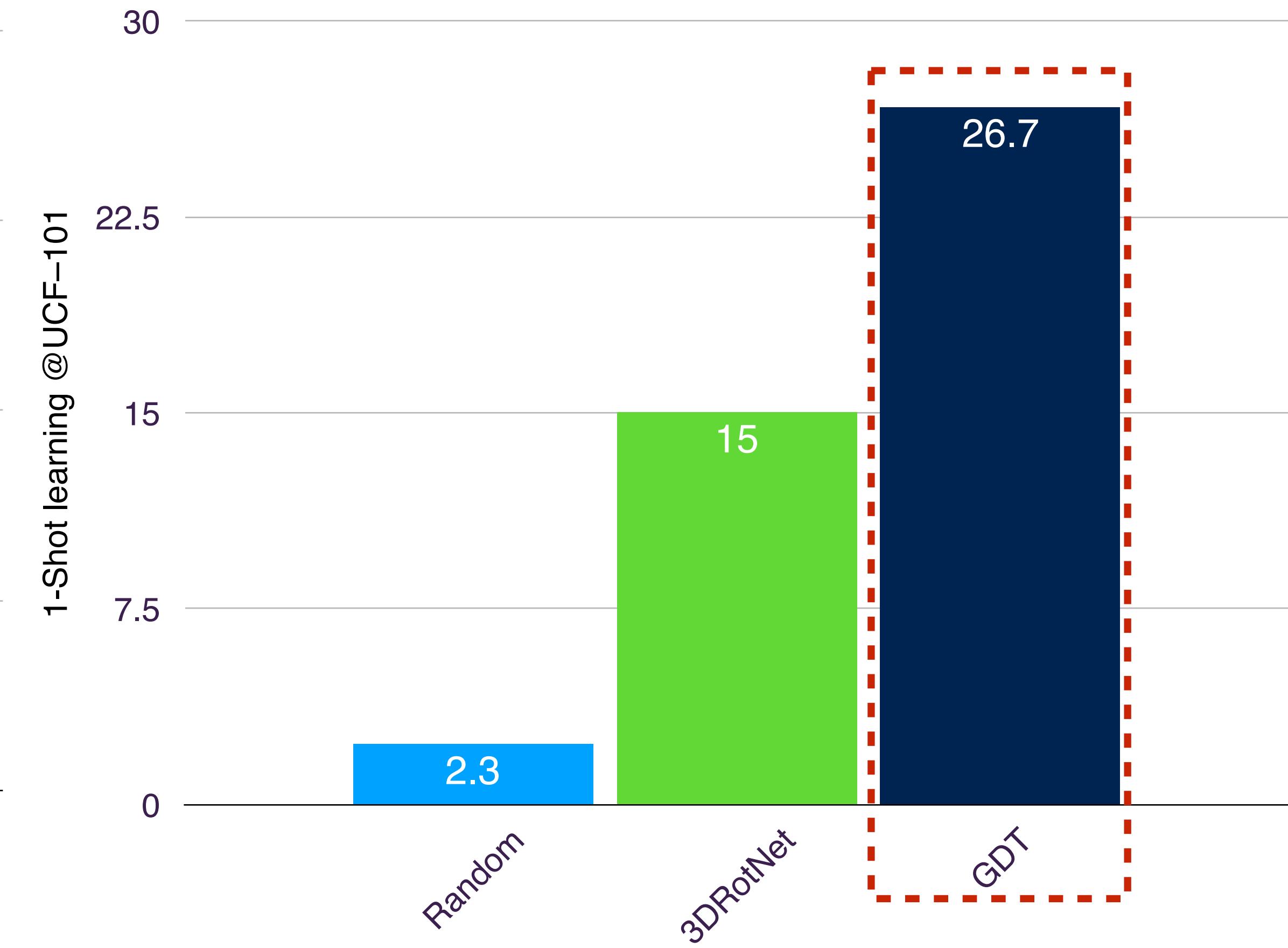
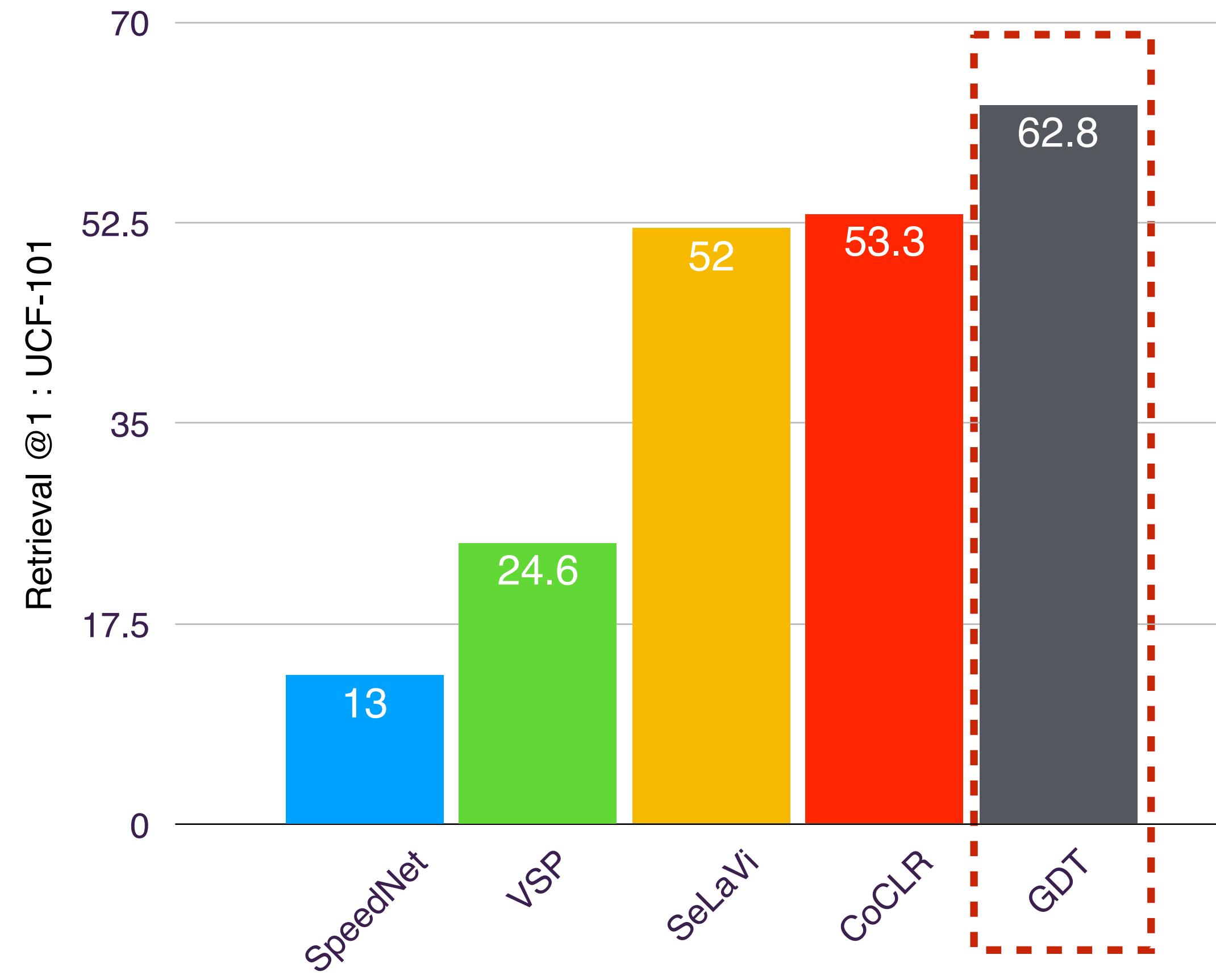


...cross-modally.

Gains from hypotheses



SOTA video action retrieval and few-shot learning results



SOTA finetuning video-action recognition results

