

## Introduction: Zero-Shot Action Recognition

- **Task:** classifying actions without any training data (**unseen target classes**)
- **How?** By linking visual and semantic features through **seen source classes**



## Contributions & Summary

### Motivation

- Recent work shows extraordinary results when using external data sources for zero-shot action recognition
- **Problem:** in a cross-dataset setup source and target categories are often not disjoint

### Contributions

- We show that external sources often **have actions excessively similar to the target classes**, strongly influencing the performance and **violating the ZSL premise**
- We propose an **evaluation procedure** that enables fair use of external data for zero-shot action recognition
- Side-contribution: we propose the **hybrid evaluation regime**, which uses the available training data of the source domain **and** the large-scale external datasets

## Fair transfer of foreign categories

### Evaluation regimes for ZSL

**Intra-dataset: same origin of training- and test data**

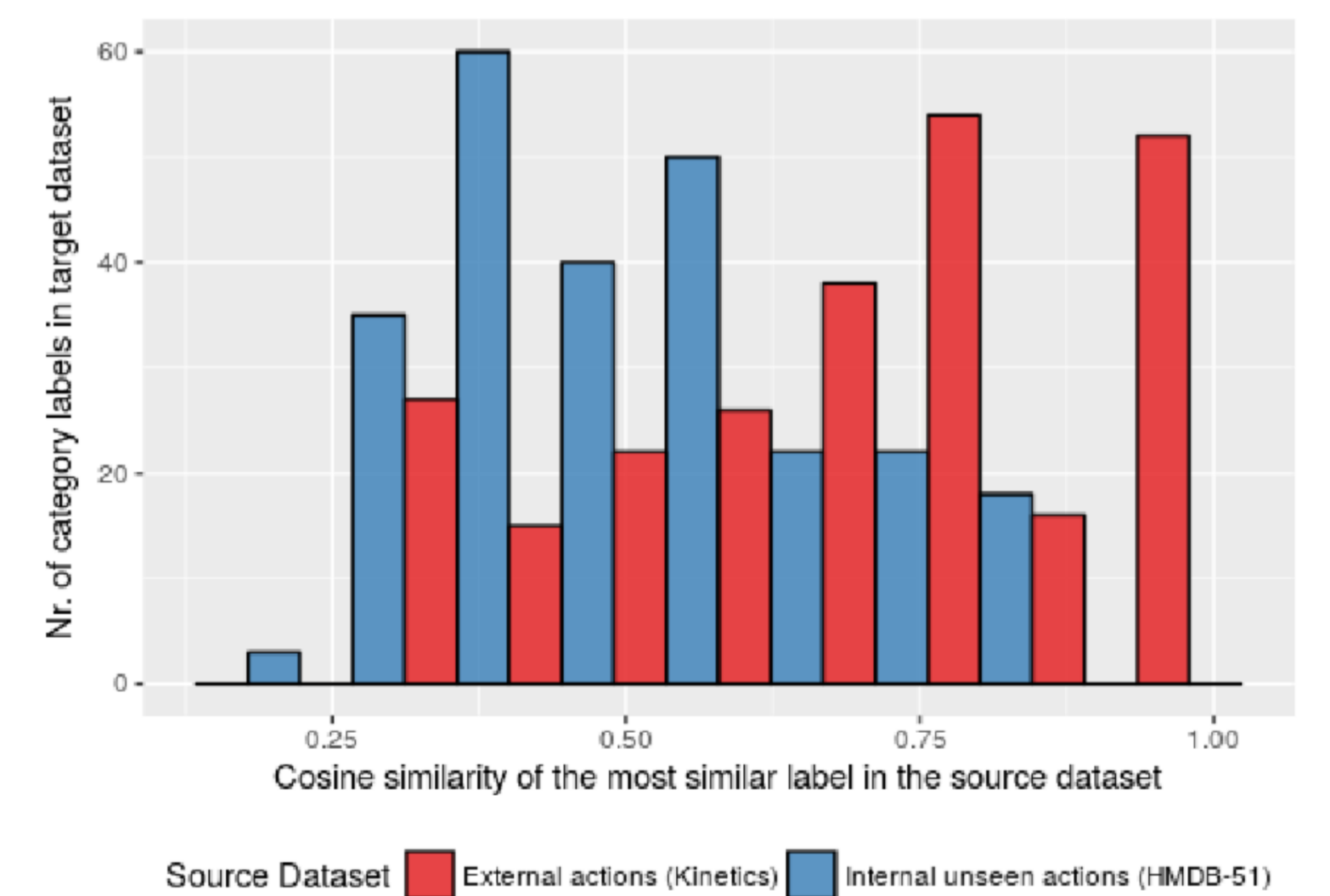
Approach	Description	Standard Approaches
Supervised Action Recognition (AR)	Classifying the already known categories $T \subset S$	standard approaches
Zero-Shot AR Intra-dataset	Source from the <b>same domain</b> : $S = S_{native}$ $T \cap S_{native} = \emptyset \rightarrow$ <b>ZSL premise satisfied</b> 😊	

**Cross-dataset: utilize large-scale external data sources**

Approach	Description	Novel Approaches
Zero-Shot AR Cross-dataset (Zhu et. al, 2018)	Source from a different domain: $S = S_{ext}$ Boost in accuracy $T \cap S_{ext} \neq \emptyset$ . <b>ZSL premise not given by default</b> ☹️ <b>Our corrective protocol eliminates too unfamiliar concepts <math>\rightarrow</math> ZSL premise given</b>	novel approaches
Zero-Shot AR Hybrid (ours)	Source from native and external domains: $S = S_{ext} \cup S_{native}$ Boost in accuracy, <b>lower-bounded by the intra- and cross-dataset regimes</b> Same evaluation issues as in cross-dataset <b>Our corrective protocol eliminates too unfamiliar concepts <math>\rightarrow</math> ZSL premise given</b>	

### Problem with the source-target synonyms

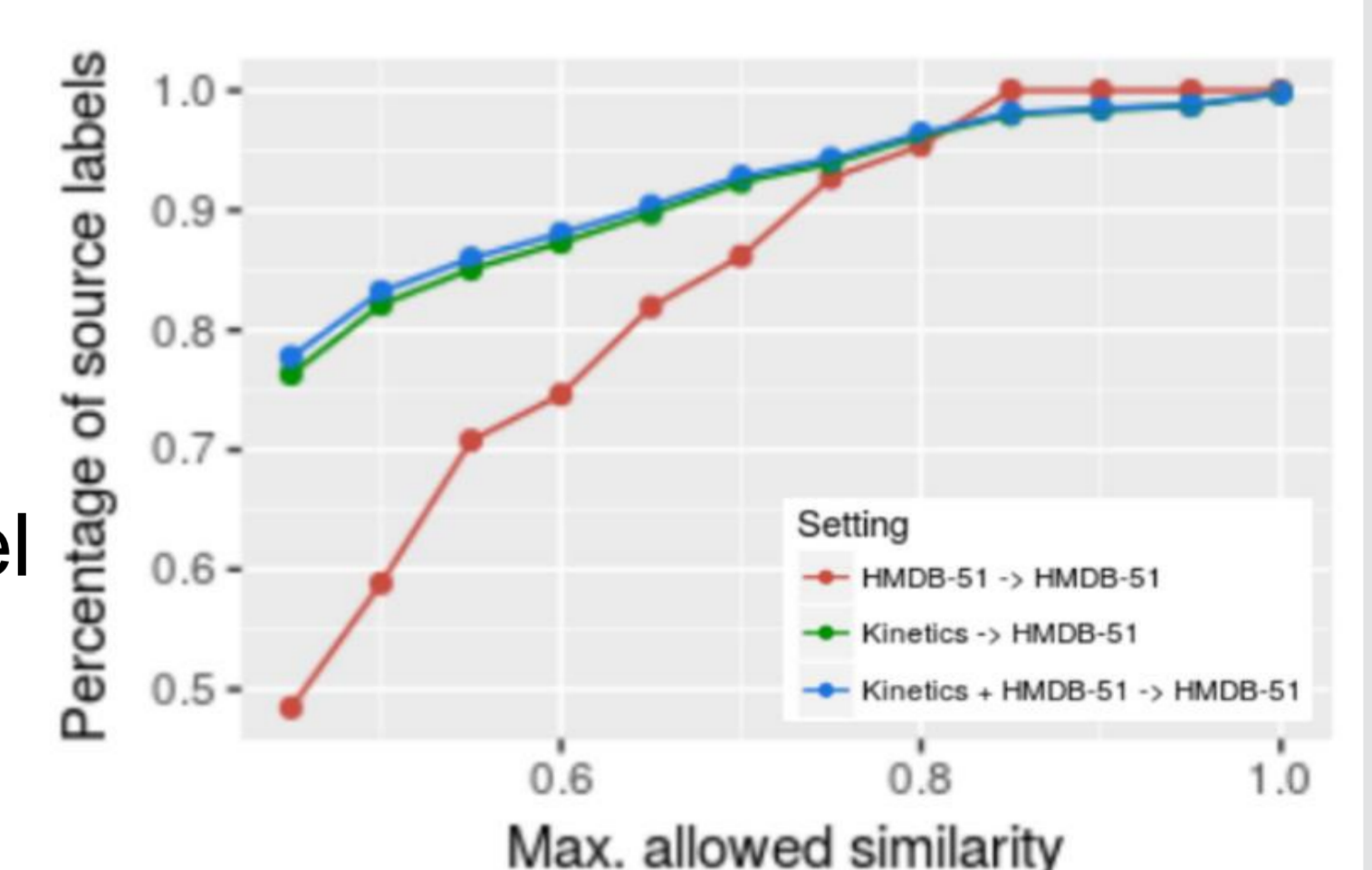
- A dataset does not contain the same action class twice
- External datasets intersect with datasets for zero-shot AR!
- Example: *brushing hair* in ActivityNet, Kinetics and HMDB51 (*brush hair*)



- Specializations: *drinking beer* vs. *drinking*  
 $\rightarrow$  **Getting rid of the direct matches is not enough!**

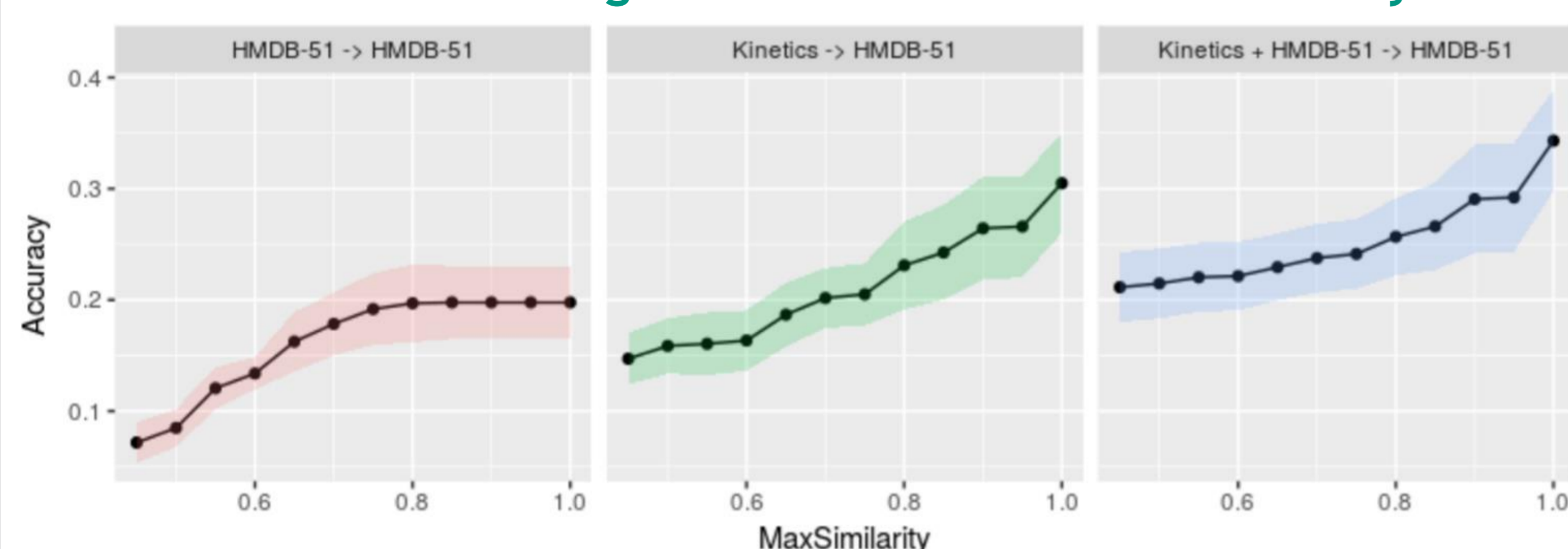
### Proposed corrective protocol for fair cross-dataset transfer

- 1) Calculate the maximum intra-dataset similarity as our threshold  $s_{th}$ :  
$$s_{th} = \max_{a_k \in S_{intra}, t_m \in T} s(\omega(a_k), \omega(t_m))$$
- 2) Purge the source category, if the label is too similar:  
$$\forall t_m \in T, s(\omega(a_k), \omega(t_m)) \leq s_{th}$$



## Experiments

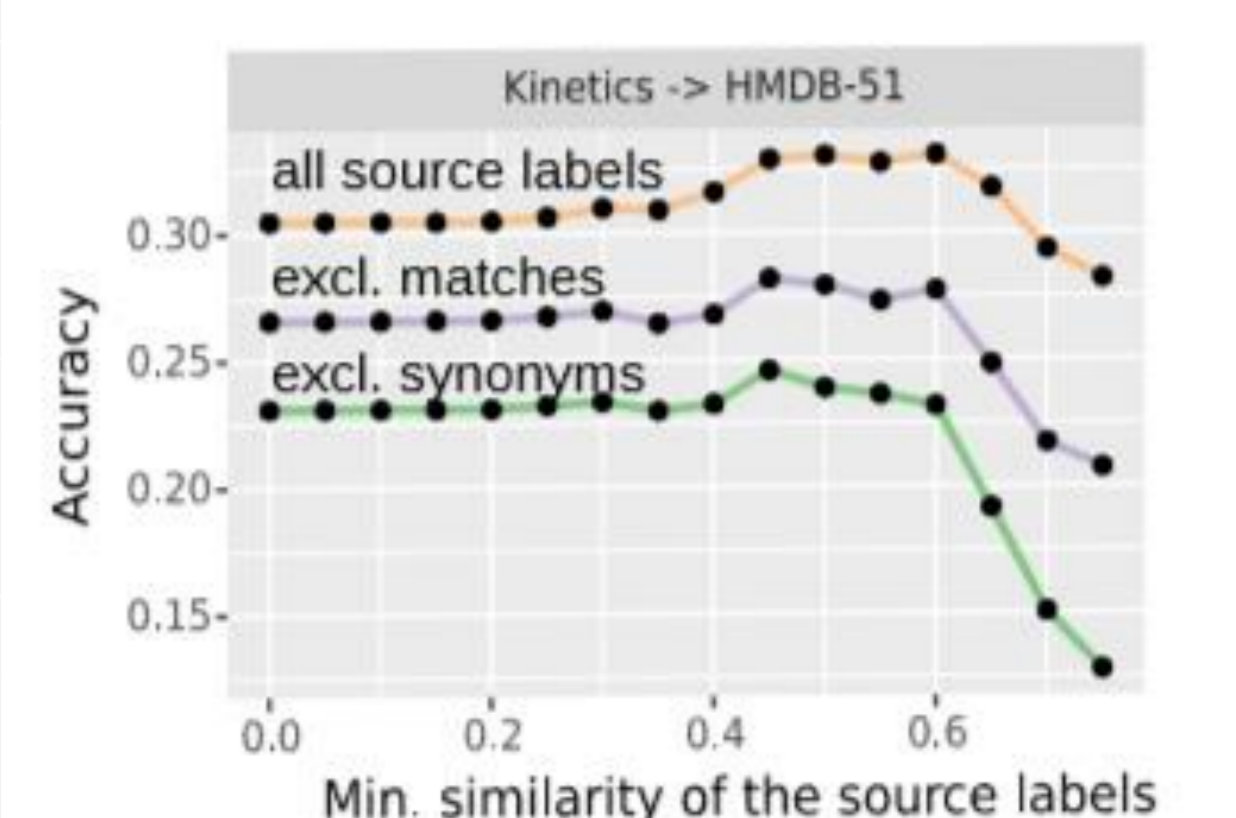
### Similar source and target classes influence the accuracy



### Setup Details

- ZSL Method: Convex Combination of Semantic Embeddings (ConSE)
- Language Model: word2vec, visual model: I3D
- Ten random splits into seen/unseen categories (26/25) for HMDB-51.
- Kinetics as external source (400 activity classes)

### Eliminating too unfamiliar concepts



An additional **lower bound** on the source similarity improves the performance