# Unsupervised Learning from Narrated Instruction Videos

**Jean-Baptiste Alayrac, Piotr Bojanowski, Nishant Agrawal, Josef Sivic, Ivan Laptev, Simon Lacoste-Julien**
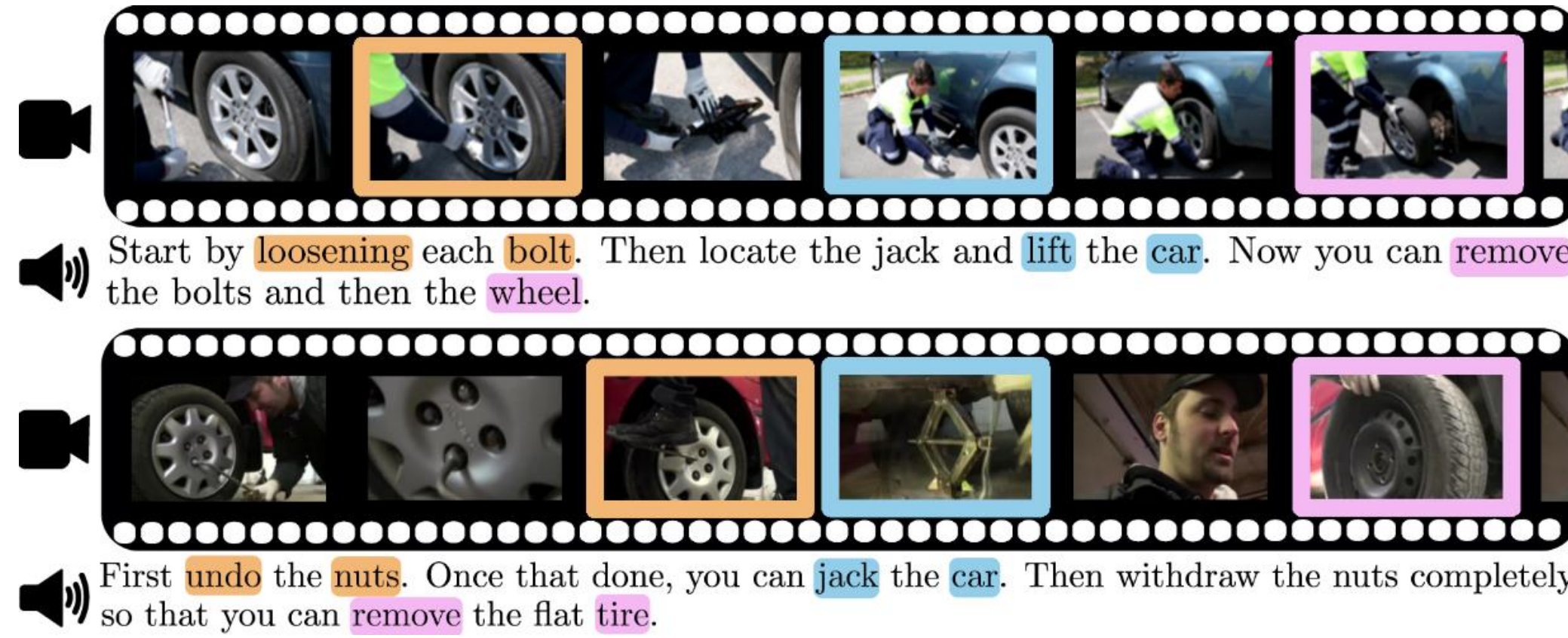
INRIA/ENS, Paris, France

## Goal and overview

**The problem:**

**Automatically** learn the **main steps** to complete a given **task** from **narrated instruction videos**.

**Input:** **A set of narrated instruction videos.**

Start by loosening each bolt. Then locate the jack and lift the car. Now you can remove the bolts and then the wheel.

First undo the nuts. Once that done, you can jack the car. Then withdraw the nuts completely so that you can remove the flat tire.

**Outputs:**

1. Loosen nuts
2. Jack the car
3. Remove the flat tire

- List of **K** (input) main steps
- **Visual and linguistic representations of the steps**
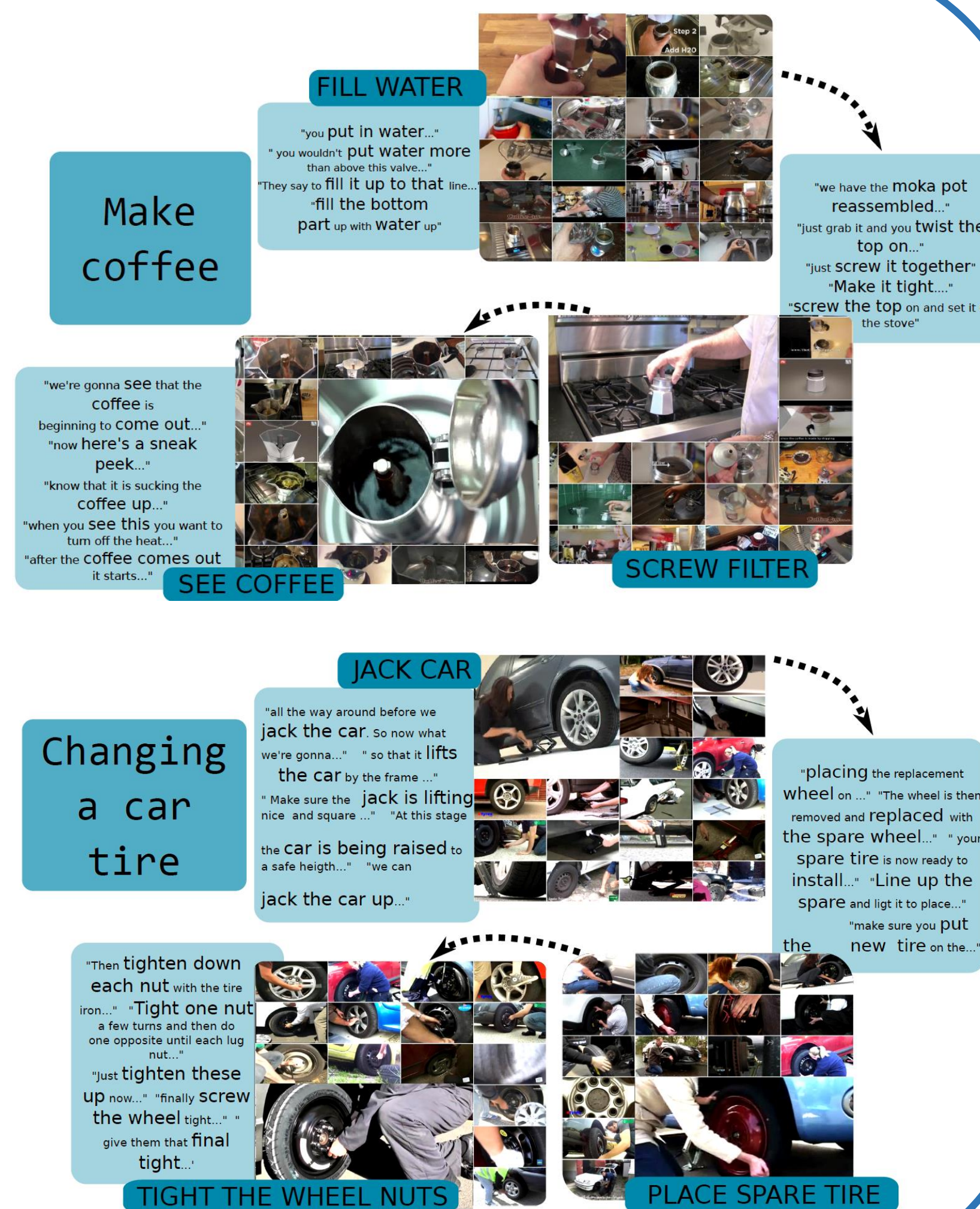- **Temporal localization of each step in the videos**

**Contributions:**

i. Collected and annotated a **new dataset** of narrated instruction videos,

ii. Developed an **unsupervised learning method** that takes advantage of the complementary nature of the text and video,

iii. Experimentally demonstrated **recovery of main steps** and their **locations in video**.

## A new dataset

- **5** tasks:
  - Changing a car tire
  - Repot a plant
  - Make coffee
  - Perform CPR
  - Jump a car battery
- **30** videos per task (total of **800,000** frames).
- Manual correction of the **ASR** transcripts.
- **Manual annotation** of 7-10 main steps and time localization in each video (only used for evaluation).

FILL WATER · Make coffee · SEE COFFEE · SCREW FILTER · JACK CAR · Changing a car tire · TIGHT THE WHEEL NUTS · PLACE SPARE TIRE

## Approach

### Assumptions

- Each task is performed by an **ordered** sequence of steps.
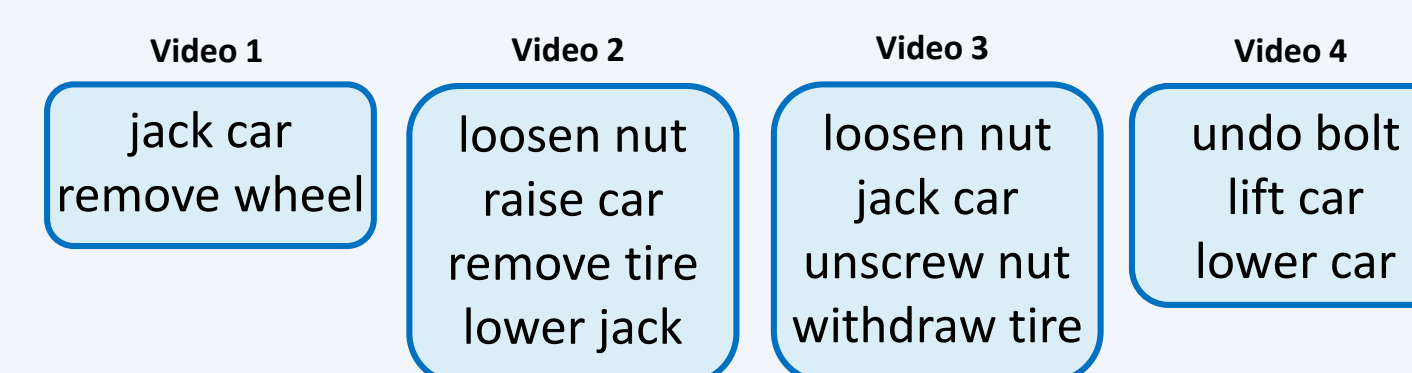- People do **what** they say (roughly) **when** they say it

### Two linked clustering steps

1. Text clustering : multiple sequence **alignment**
2. Discriminative video clustering under **text constraints**
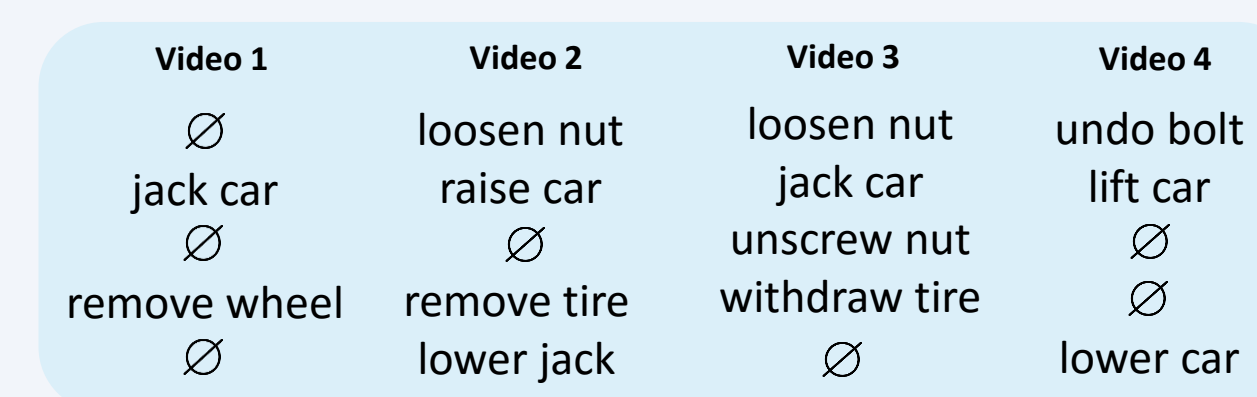
### 1st STEP : TEXT — Identify main steps

#### Multiple sequence alignment (MSA)

- Text signals are first processed into sequences of **direct object relations**: *Let's now jack the car.* → **DOBJ = (jack,car)**
- Similarity between **dobj** is obtained using **Wordnet**.

**INDIVIDUAL INPUT SEQUENCES:**

| Video 1 | Video 2 | Video 3 | Video 4 |
|---|---|---|---|
| jack car | loosen nut | loosen nut | undo bolt |
| remove wheel | raise car | jack car | lift car |
| | remove tire | unscrew nut | lower car |
| | lower jack | withdraw tire | |

**OUTPUT OF THE MSA:**

**New!** Formulate **MSA** as a **QP**, and approximately solve it with **Frank-Wolfe**

| Video 1 | Video 2 | Video 3 | Video 4 |
|---|---|---|---|
| ∅ | loosen nut | loosen nut | undo bolt |
| jack car | raise car | jack car | lift car |
| ∅ | ∅ | unscrew nut | ∅ |
| remove wheel | remove tire | withdraw tire | ∅ |
| ∅ | lower jack | ∅ | lower car |

**FINAL OUTPUT:**

| Video 1 | Video 2 | Video 3 | Video 4 | Agreement | Discovered list of steps |
|---|---|---|---|---|---|
| ∅ | loosen nut | loosen nut | undo bolt | 3 | 1) Loosen nut |
| jack car | raise car | jack car | lift car | 4 | 2) Jack car |
| ∅ | ∅ | unscrew nut | ∅ | 1 | 3) Remove wheel |
| remove wheel | remove tire | withdraw tire | ∅ | 3 | |
| ∅ | lower jack | ∅ | lower car | 2 | |

*Text assignment [SxK] matrix*

### 2nd STEP : VIDEO — Temporally localize each steps

#### Discriminative clustering under text constraints

**CONSTRAINED OPTIMIZATION PROBLEM:**

$$\underset{Z}{\text{minimize}} \quad h(Z) \quad \text{s.t.} \quad Z \in \mathcal{Z}, \quad AZ \geq R.$$

- Discovered temporal localization [TxK] matrix
- ordered script
- weak textual constraints
- Subtitle Alignment [SxT] matrix

where $h(Z)$ is a **discriminative clustering** cost [1]:

$$h(Z) = \underset{W \in \mathbb{R}^{K \times d}}{\min} \frac{1}{2T}\|Z - XW\|_F^2 + \frac{\lambda}{2}\|W\|_F^2.$$

- Representation of video (IDTF,CNN) [Txd] matrix
- Discriminative loss on data
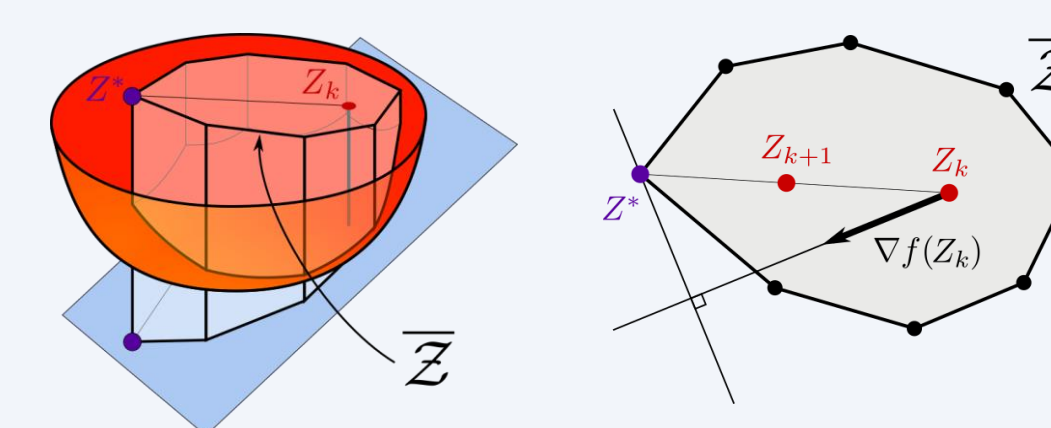- Regularizer
- Linear action classifier [dxK] matrix

T = total number of time intervals
S = total number of text tokens
K = number of steps
d = dimension of the features

**OPTIMIZATION METHOD [2,3]:**

- Optimize convex relaxation using Frank-Wolfe
- Use **DP** as the linear oracle
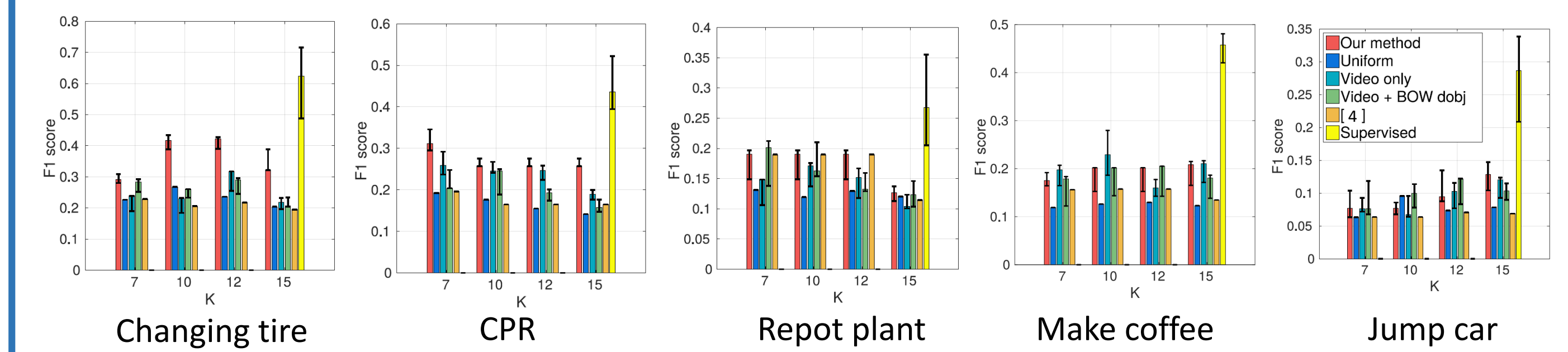- Cost classifier based **rounding**

## Experiments

### Script discovery results

| Changing a tire | | Performing CPR | | Repot a plant | | Make coffee | | Jump car | |
|---|---|---|---|---|---|---|---|---|---|
| GT (11) | K ≤ 10 | GT (7) | K ≤ 10 | GT (7) | K ≤ 10 | GT (10) | K ≤ 10 | GT (12) | K ≤ 10 |
| get tools out | | open airway | open airway | take plant | remove plant | add coffee | put coffee | connect red A | connect cable |
| start loose | | check pulse | put hand | put soil | use soil | | fill chamber | | charge battery |
| | | | tilt head | loosen roots | loosen soil | | fill water | connect red B | connect end |
| jack car | | | lift chin | place plant | place plant | fill water | put filter | start car A | start car |
| remove wheel | | give breath | give breath | add top | add soil | screw filter | see steam | remove cable A | remove cable |
| unscrew wheel | remove nut | do compressions | do compression | water plant | water plant | add top | take minutes | remove cable B | disconnect cable |
| remove wheel | take wheel | take breath | open airway | | | put stove | make coffee | | |
| put wheel | take tire | | put nut | | | | make cup | | |
| screw wheel | put nut | | start compression | | | see coffee | | | |
| lower car | lower jack | | do compression | | | pour coffee | | | |
| tight wheel | tighten nut | | give breath | | | | | | |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Precision | 0.9 | Precision | 0.4 | Precision | 1 | Precision | 0.67 | Precision | 0.83 |
| Recall | 0.9 | Recall | 0.57 | Recall | 0.86 | Recall | 0.6 | Recall | 0.42 |

### Localizing instruction steps in video

Changing tire · CPR · Repot plant · Make coffee · Jump car

Our method · Uniform · Video only · Video + BOW dobj · Supervised

### Qualitative results

CHANGING CAR TIRE · UNSCREW JACK CAR WHEEL · LOWER CAR · REPOTING A PLANT · PLACE PLANT · ADD TOP · LOOSEN ROOT

PERFORMING CPR · OPEN AIRWAY · GIVE BREATH · GIVE COMPR. · MAKING COFFEE · SCREW TOP · FILL WATER · SEE COFFEE

JUMPING CARS · CONN. RED A · CONN. GROUND · DISC. RED B

## References

[1] Bach and Harchaoui. DIFFRAC: A discriminative and flexible framework for clustering. In *NIPS, 2007*.

[2] Bojanowski et al. Weakly supervised action labeling in videos under ordering constraints. In *ECCV, 2014*.

[3] Bojanowski et al. Weakly-Supervised Alignment of Video With Text. In *ECCV, 2015*.

[4] Malmaud et al. What's cookin'? Interpreting cooking videos using text, speech and vision. In *NACL, 2015*.

**Check out our project webpage for code/data!**