# Detailed Human Action Understanding from Unlabeled Videos

Chen Sun

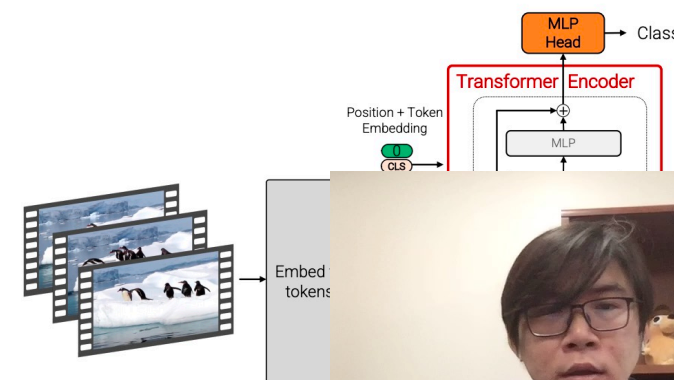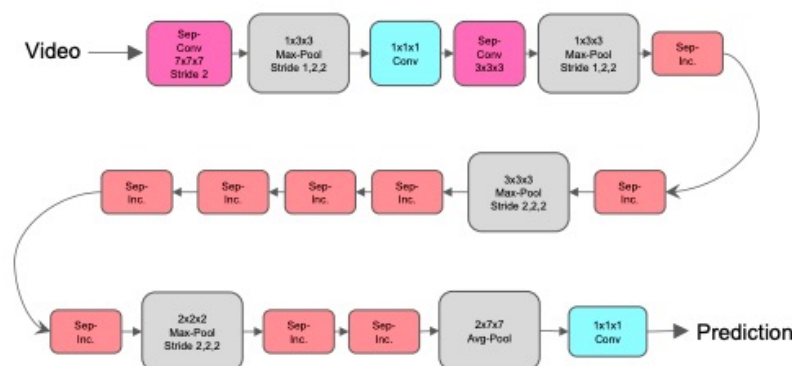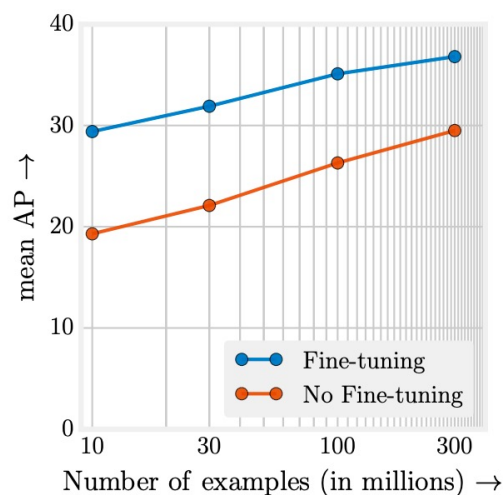BROWN        Google Research
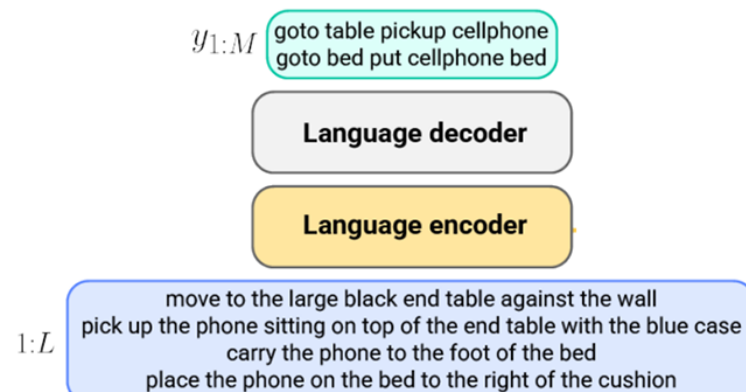
# My Research at Google: Large-scale Visual Understanding



riding mountain bike:0.996515
crossing river:0.000579438
riding a bike:0.000265659



Left: Stand, Watch; Middle: Stand, Play instrument; Right: Sit, Play instrument

# My Research at Brown: Structured Video Understanding



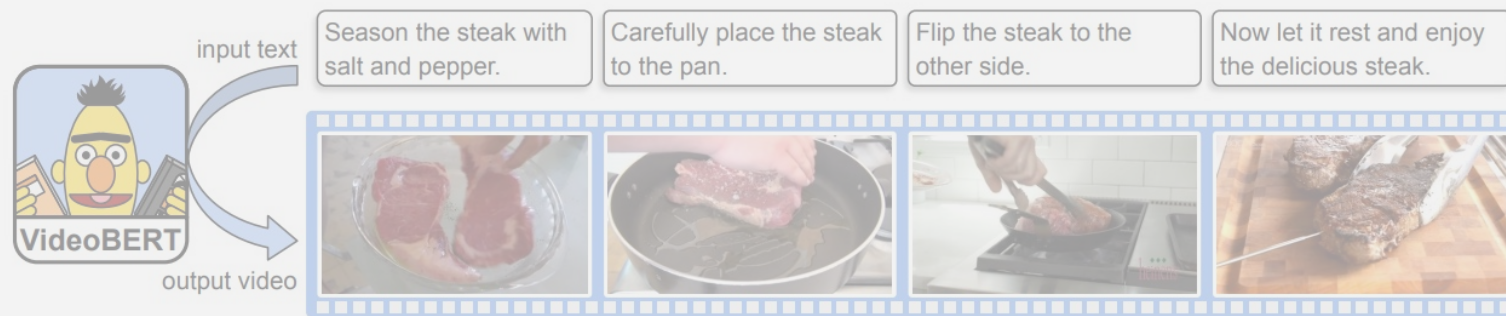| Season the steak with salt and pepper. | Carefully place the steak to the pan. | Flip the steak to the other side. | Now let it rest and enjoy the delicious steak. |

Goal: "put two vases on a cabinet"

t=0 — "Walk forwards and then turn right. Pick up the vase from the fireplace."

t=12 — "Turn right and then left."

t=30 — "Put the vase on the cabinet."

t=31 — "Go to the right of the fire-place. Pick up another vase."

t=40 — "Walk back to where you were standing previously with the second vase."

t=53 — "Put the second vase on the same cabinet."

$y_{1:M}$ — goto table pickup cellphone goto bed put cellphone bed

Language decoder

Language encoder

$1:L$ — move to the large black end table against the wall
pick up the phone sitting on top of the end table with the blue case
carry the phone to the foot of the bed
place the phone on the bed to the right of the cushion

We are hiring PhD students!

# What can we learn from videos?



A frame from the Atomic Visual Actions (AVA) dataset

# What can we learn from videos?



A frame from the Atomic Visual Actions (AVA) dataset

**Object detection**:
*Person, silverware, food*
**Action detection**:
*Sit, eat, talk*
**Human-object interaction**:
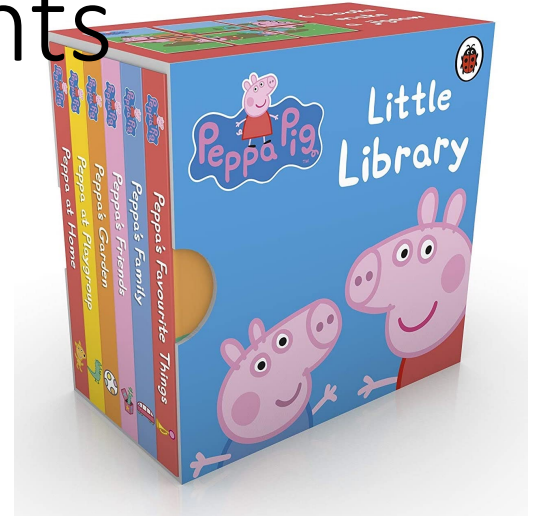*Person hold fork / eat food*
**Near-future prediction**:
*Stand*

# Encyclopedia of Multimedia Contents

Place the ingredients onto a bowl of hot steamed rice.

# What **else** can we learn from videos?

How to Turn  into:



OMELET

BAKED

POACHED

SCRAMBLED

HARD BOILED

FRIED

OVER EAS

# What **else** can we learn from videos?



**Transfer what has been learned from passive observations**

# A RoadMap Towards Video Understanding



**Object Representation** → **Scene Representation** → **Short Video Representation** → **Multimodal Representation** → **Temporal Dynamics** → **Video Understanding**

*Hard to label all of them*
*Need to learn from unlabeled videos!!*

# Scene-level Contrastive Learning

View 1: Augmented image

View 2: Augmented image

Similar

Different

DistInst
CPC
CMC
SimCLR
MoCo

...

# Contrastive Learning for Videos



Qian and Meng et al., Spatiotemporal Contrastive Video Representation Learning, CV

# What should consist positive pairs?

For images:
Preserve objects

For videos:
?

# Natural views introduce undesired invariances

View 2: $v^{t+t'}$

View 1: $v^t$

Invariant to "noise"

Representation space

# Natural views introduce undesired invariances

View 2: $v^{t+t'}$



View 1: $v^t$

Invariant to "noise"

Representation space

**Signal:**
Color, local flow

**"Noise":**
Shape deformation

Loses temporal info

# Solution 1: Construct many pairs of views



May not scale well

Xiao et al., What Should Not Be Contrastive in Contrastive Learning, ICLR 202

# Solution 2: Equivariant representations

Not necessary for many tasks

Jayaraman and Grauman, Learning image representations tied to ego-motion, ICCV

# Our solution: Simply encode the augmentations

View 2: $v^{t+t'}$



View 1: $v^t$



Representation space

# Our solution: Simply encode the augmentations

View 2: $v^{t+t'}$



View 1: $v^t$



Rewind(*t'*)

Representation space

Learn an implicit "prediction" model of t'

**Shared** and **predictable** information can be preserved: color, shape, etc.

Unpredictable is still "noise" ca...

**Special cases:** view-invariant coding, view-predictive coding

# Composable AugmenTation Encoding (CATE)



Projection head is now a Transfor[mer]
encodes a sequence of augmenta[tions]

Sun, Nagrani, Tian, and Schmid, Composable augmentation encoding for video representation l[earning]

# The Something-Something Dataset



Classes

| | |
|---|---|
| Putting something on a surface | 4,081 |
| Moving something up | 3,750 |
| Covering something with something | 3,530 |
| Pushing something from left to right | 3,442 |
| Moving something down | 3,242 |
| Pushing something from right to left | 3,195 |
| Uncovering something | 3,004 |
| Taking one of many similar things on the table | 2,969 |

Fine-grained actions that rely on
the arrow of ti

# Augmentation encoding is helpful

| Encoded | $\tau$ | Dropout | Top-1 Acc. | Top-5 Acc. |
|---------|--------|---------|------------|------------|
| No | - | - | 26.5 | 55.9 |
| Crop | $\delta_{x,y}$ | ✗ | 27.2 | 56.7 |
| Crop | $\delta_{x,y}$ | ✓ | 28.1 | 58.0 |
| Time | $\mathrm{sgn}(\delta_t)$ | ✗ | 28.1 | 57.9 |
| Time | $\delta_t$ | ✗ | 31.3 | 62.4 |
| Time | $\delta_t$ | ✓ | 31.2 | 61.4 |

| Encode Time | $\tau$ | Time Offset Acc. |
|-------------|--------|------------------|
| ✗ | - | 5.7 |
| ✓ | $\mathrm{sgn}(\delta_t)$ | 65.7 |
| ✓ | $\delta_t$ | **99.9** |

Table 5: **Time Shift Classification on SSv1**. Encoding time significantly helps on this proxy task, validating the intuition that our model retains useful time information.

# Augmentation encoding is composable

| Enc. Crop | Enc. Time | Top-1 Acc. | Top-5 Acc. |
|:---------:|:---------:|:----------:|:----------:|
| ✗ | ✗ | 26.5 | 55.9 |
| ✓ | ✗ | 28.1 | 58.0 |
| ✗ | ✓ | 31.2 | 61.4 |
| ✓ | ✓ | **32.2** | **62.4** |

Table 2: **Composing spatial (crop) and temporal encodings** for Something-Something v1. Each individual encoding outperforms the no encoding baseline (SimCLR++). Composing them together yields the best performance.

# Per-class comparison (temporal aug.)

Arrow of time barely matters:

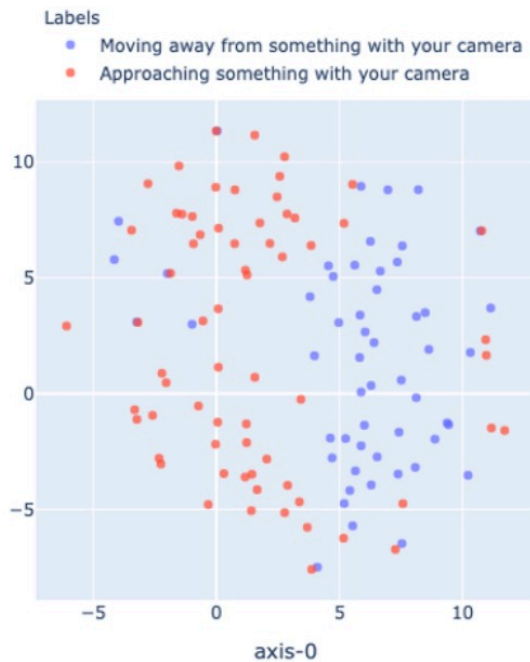| Label | $\Delta$AP |
|---|---|
| Lifting something up completely, then letting it drop down | 21.0 |
| Pulling two ends of something so that it gets stretched | 19.8 |
| Moving something and something closer to each other | 18.5 |
| Taking one of many similar things on the table | 17.2 |
| Pushing something so that it almost falls off but doesn't | 16.7 |
| Poking something so lightly that it doesn't move | -4.6 |
| Pretending to pour something out of something | -5.4 |
| Poking a stack of something without the stack collapsing | -5.5 |
| Pretending to spread air onto something | -7.8 |

Table 4: Classes that benefit the most and the least with **time encoding** on SSv1. We sort the classes by their differences on Average Precision.
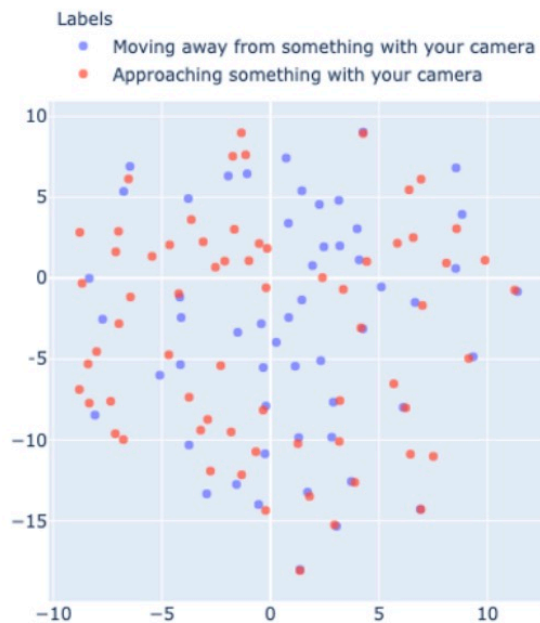
# t-SNE



CATE

No encoding

# Comparison on other benchmarks

| Method | top 1 | top 5 | top 10 | top 20 | top 50 |
|---|---|---|---|---|---|
| OPN [32] | 19.9 | 28.7 | 34.0 | 40.6 | 51.6 |
| SpeedNet [5] | 13.0 | 28.1 | 37.5 | 49.5 | 65.0 |
| VCP [34] | 19.9 | 33.7 | 42.0 | 50.5 | 64.4 |
| Temporal SSL [25] | 26.1 | 48.5 | 59.1 | 69.6 | 82.8 |
| MemDPC$^\dagger$ [18] | 40.2 | 63.2 | 71.9 | 78.6 | - |
| CATE | **54.9** | **68.3** | **75.1** | **82.3** | **89.9** |

Table A6: Nearest neighbor retrieval evaluation on UCF-101 split 1. $\dagger$: with Flow

| Method | top 1 | top 5 | top 10 | top 20 | top 50 |
|---|---|---|---|---|---|
| VCP [34] | 6.7 | 21.3 | 32.7 | 49.2 | 73.3 |
| MemDPC$^\dagger$ [18] | 15.6 | 37.6 | 52.0 | 65.3 | - |
| CATE | **33.0** | **56.8** | **69.4** | **82.1** | **92.8** |

Table A7: Nearest neighbor retrieval evaluation on HMDB-51 split 1. $\dagger$: with Flow

# Checkpoints are released!

https://github.com/google-research/google-research/tree/master/cate
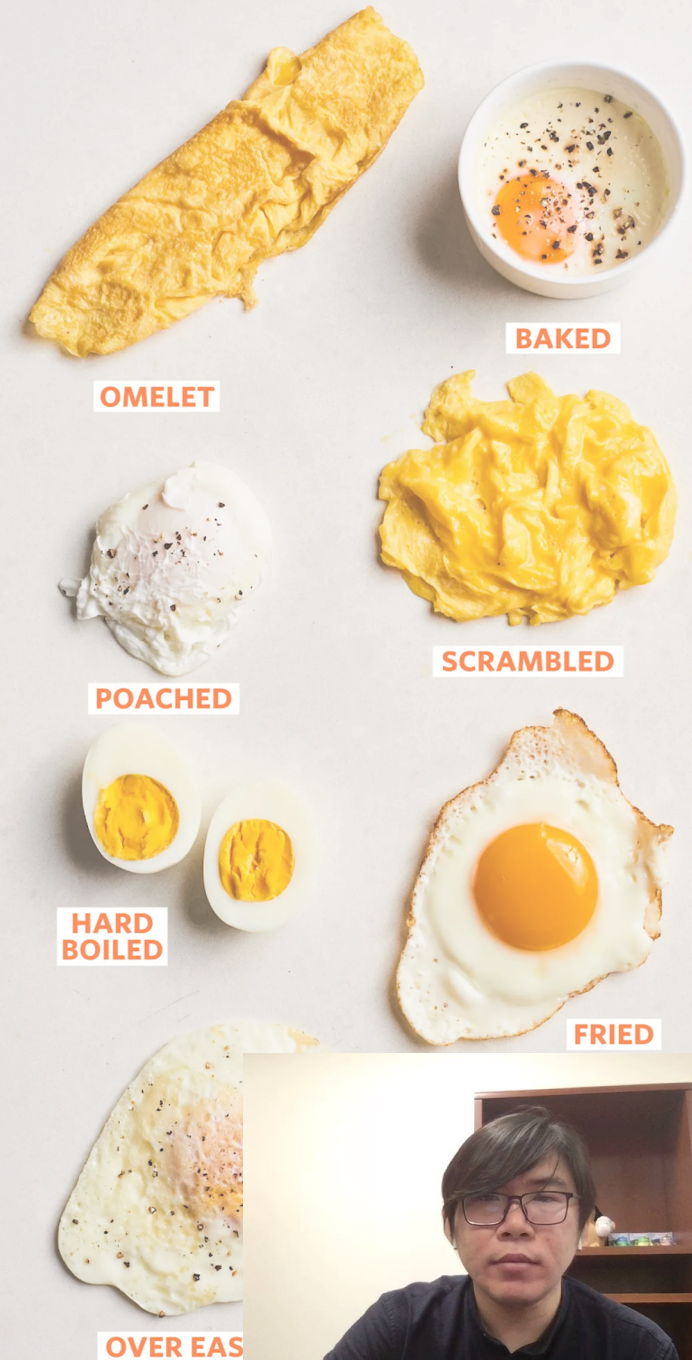
# The egg problem

$$f\left( \text{🥚}, \text{boil} \right) = \text{🥚🥚}$$

A more compact representation for videos:
**Actions as object state transitions**
(Action recognition, object tracking, ...,
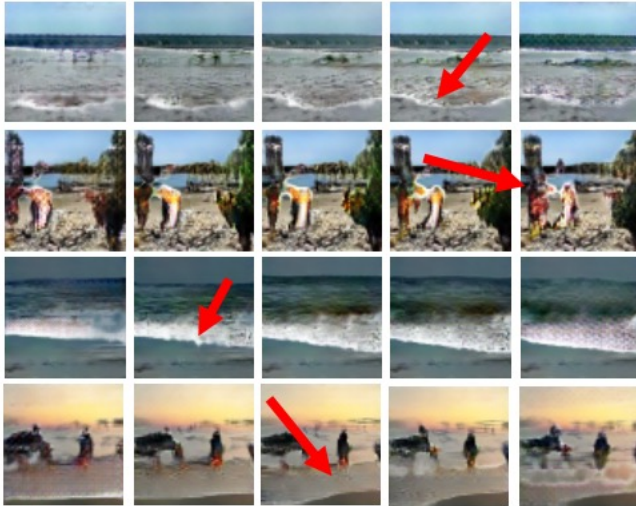Visual Commonsense)

# But why?

- Towards Long Video Understanding
    - Only use "key moments"
    - Video summarization
- Structured Representation
    - Objects
    - Their state transitions over time (visual dynamics)
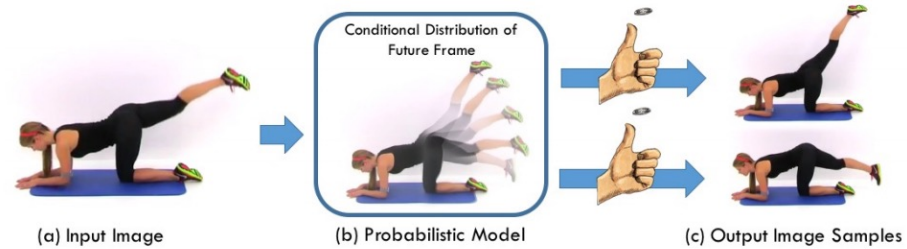- Modeling temporal dynamics is itself important

# How to predict the future?

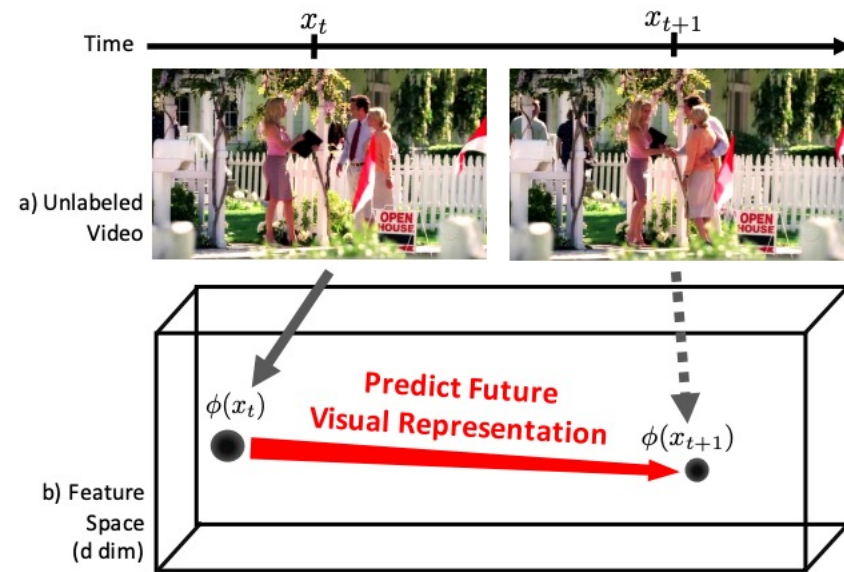Generate images…



Vondrick et al., 2016



(a) Input Image   (b) Probabilistic Model   (c) Output Image Samples
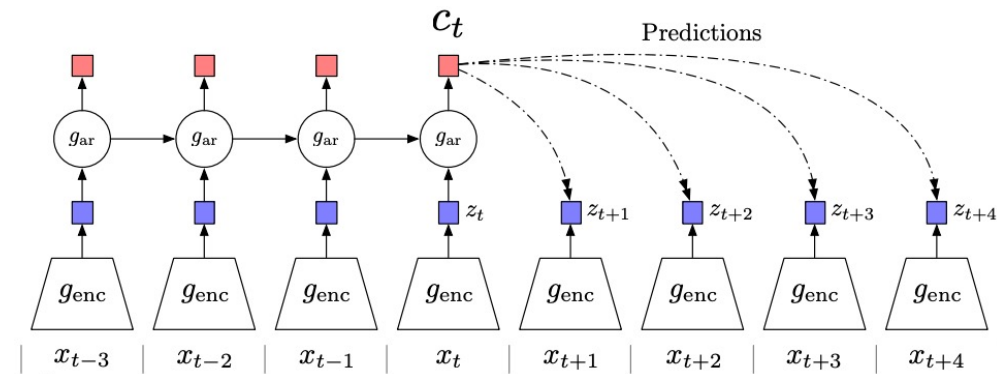
Xue et al., 2016

# How to predict the future?

Generate representations…
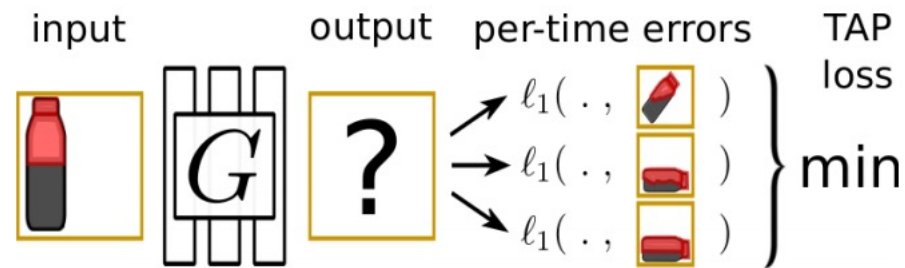


Vondrick et al., 2015

van den Oord et al., 2018

# Problem solved?

Not quite…

Predict at fixed offset into future = deal with high uncertainty!

Could let network output most predictable moment in near future
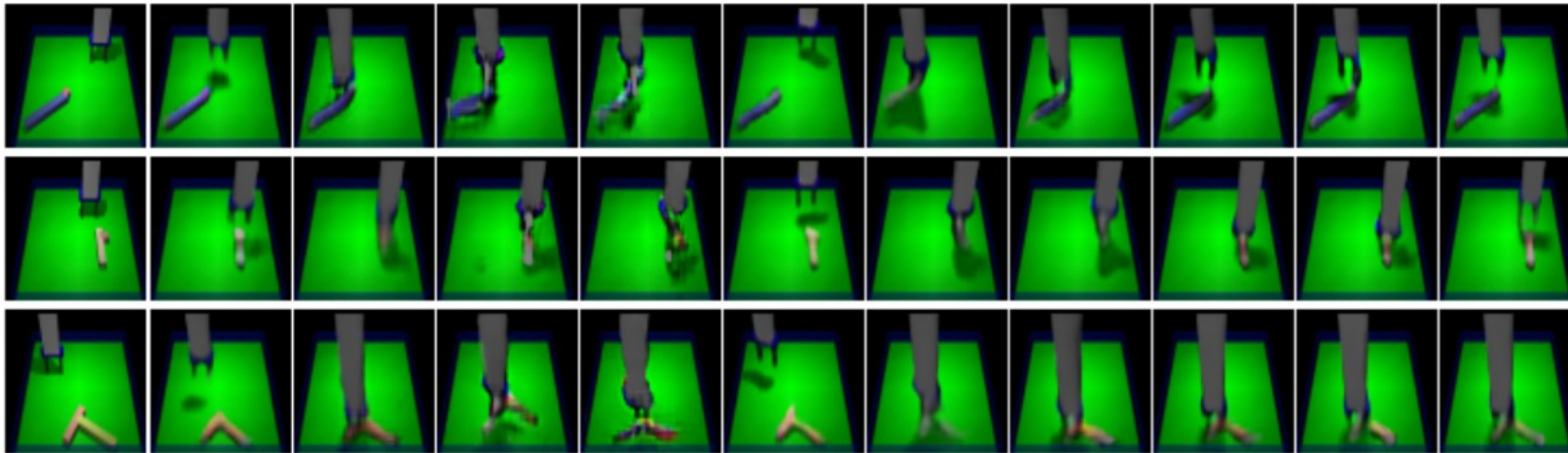


Jayaraman et al., 2018

# Okay, problem solved now?

Not quite…

Very short-term prediction – a few seconds into future at most

Limited to simple, low-level visual data



Jayaraman et al., 2018

# The ideal future prediction

Dynamic, rather than at a fixed offset into the future

High-level, e.g., mixing eggs and flour → rolling out dough

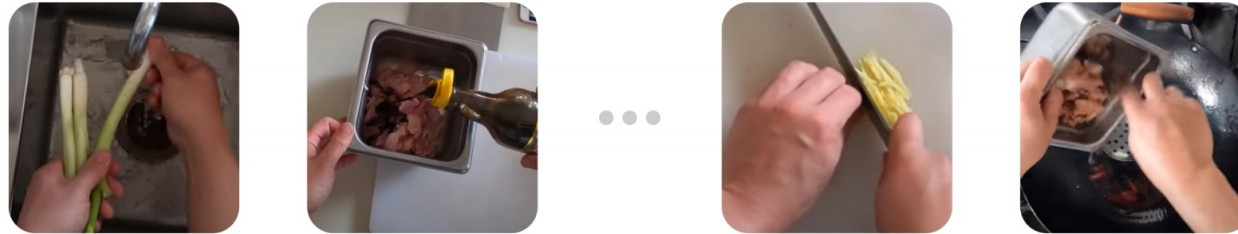**Unsupervised, to take advantage of large unlabeled datasets**

(a) Time = **t**



"go ahead and pour the cream in"

# Better future predictions



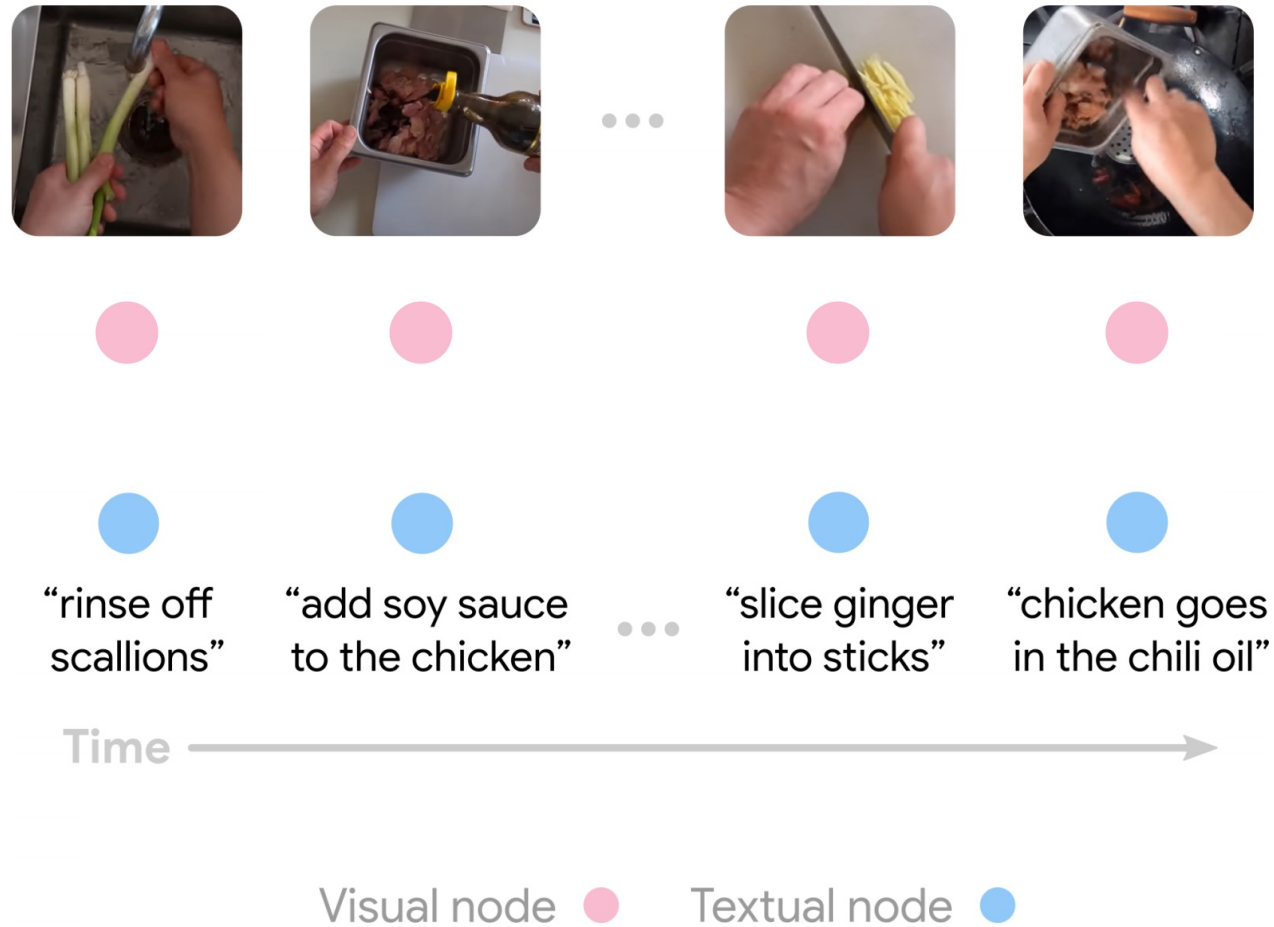"rinse off scallions"    "add soy sauce to the chicken"    ...    "slice ginger into sticks"    "chicken goes in the chili oil"

Time

Epstein, Wu, Schmid, and Sun, Learning Temporal Dynamics from Cycles in Narrated Vide

# Better future predictions



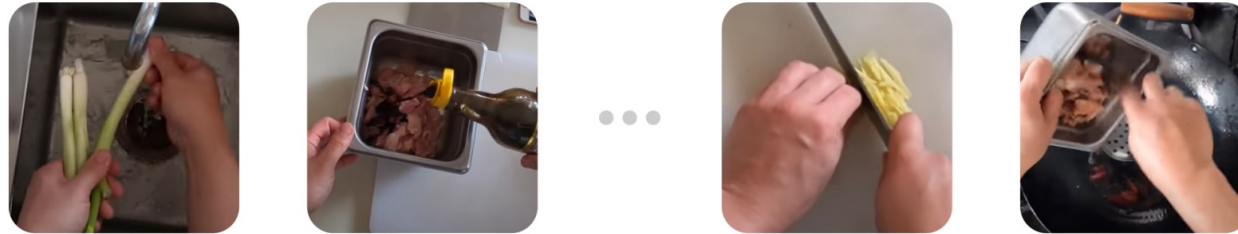"rinse off scallions"  "add soy sauce to the chicken"  ...  "slice ginger into sticks"  "chicken goes in the chili oil"

Time →

Visual node ●    Textual node ●

Epstein, Wu, Schmid, and Sun, Learning Temporal Dynamics from Cycles in Narrated Vide

# Cycling through video



"rinse off scallions"    "add soy sauce to the chicken"  ...  "slice ginger into sticks"    "chicken goes in the chili oil"
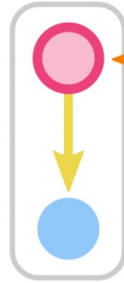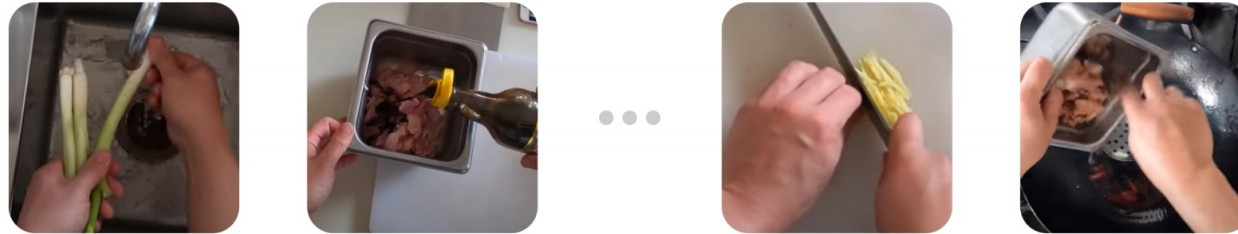
Time →

Start node ◯   Visual node ●   Textual node ●

Epstein, Wu, Schmid, and Sun, Learning Temporal Dynamics from Cycles in Narrated Vide

# Cycling through video



"rinse off scallions"   "add soy sauce to the chicken"   ...   "slice ginger into sticks"   "chicken goes in the chili oil"
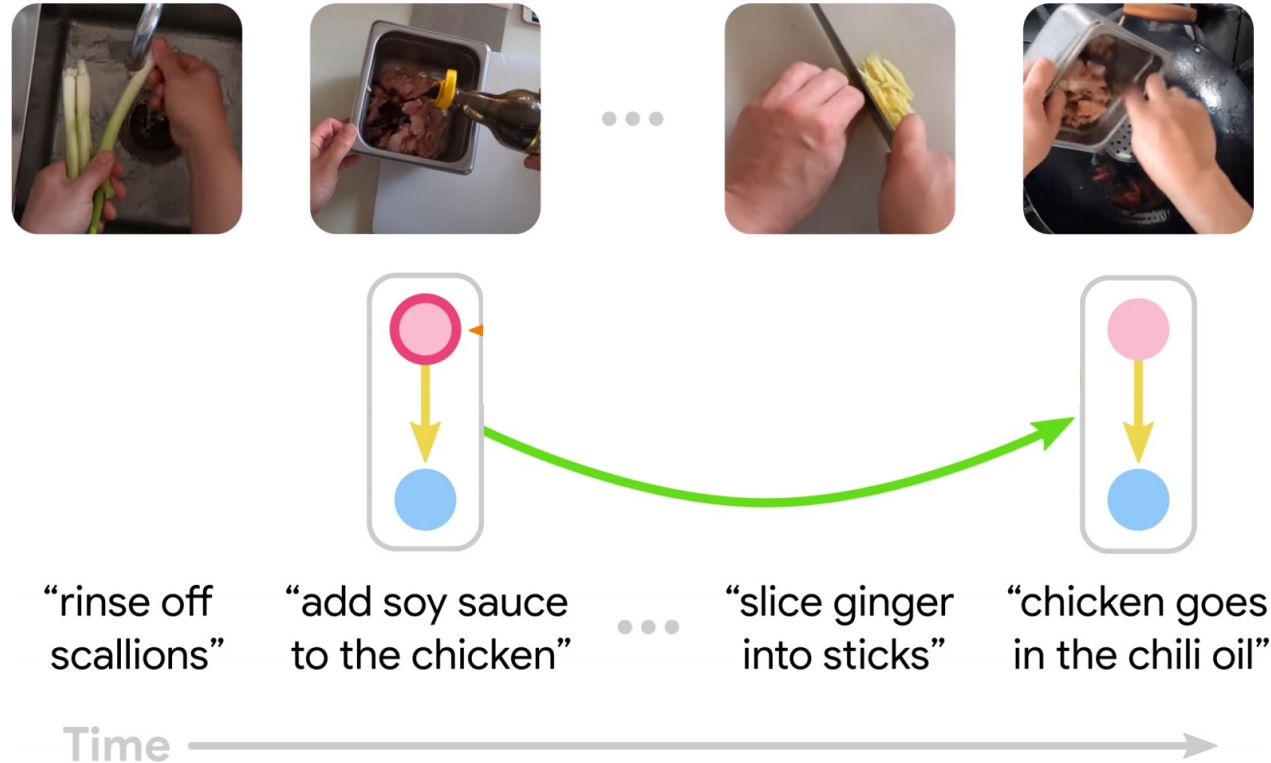
Time →

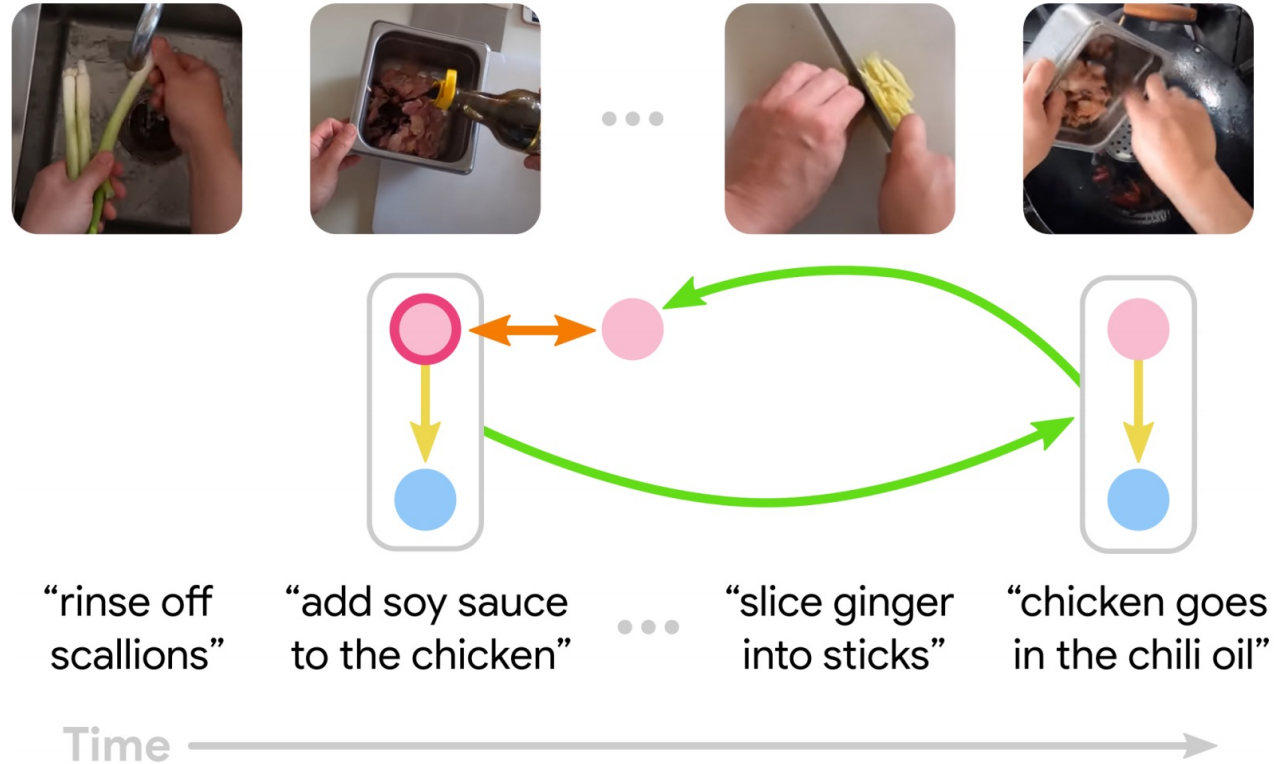Start node ○   Visual node ●   Textual node ●

Cross modal →

# Cycling through video



"rinse off scallions"      "add soy sauce to the chicken"      ...      "slice ginger into sticks"      "chicken goes in the chili oil"

Time

Start node ○      Visual node ●      Textual node ●

Cross modal →      Temporal →

# Cycling through video



"rinse off scallions"    "add soy sauce to the chicken"    ...    "slice ginger into sticks"    "chicken goes in the chili oil"

Time ⟶

Start node ○    Visual node ●    Textual node ●

Cross modal ⟶    Temporal ⟶    Loss ⟷

# Cycling through video - intuition



(a) Time = **t**
"go ahead and pour the cream in"

(b) Time = **t+1**
"go ahead and pour the cream in"

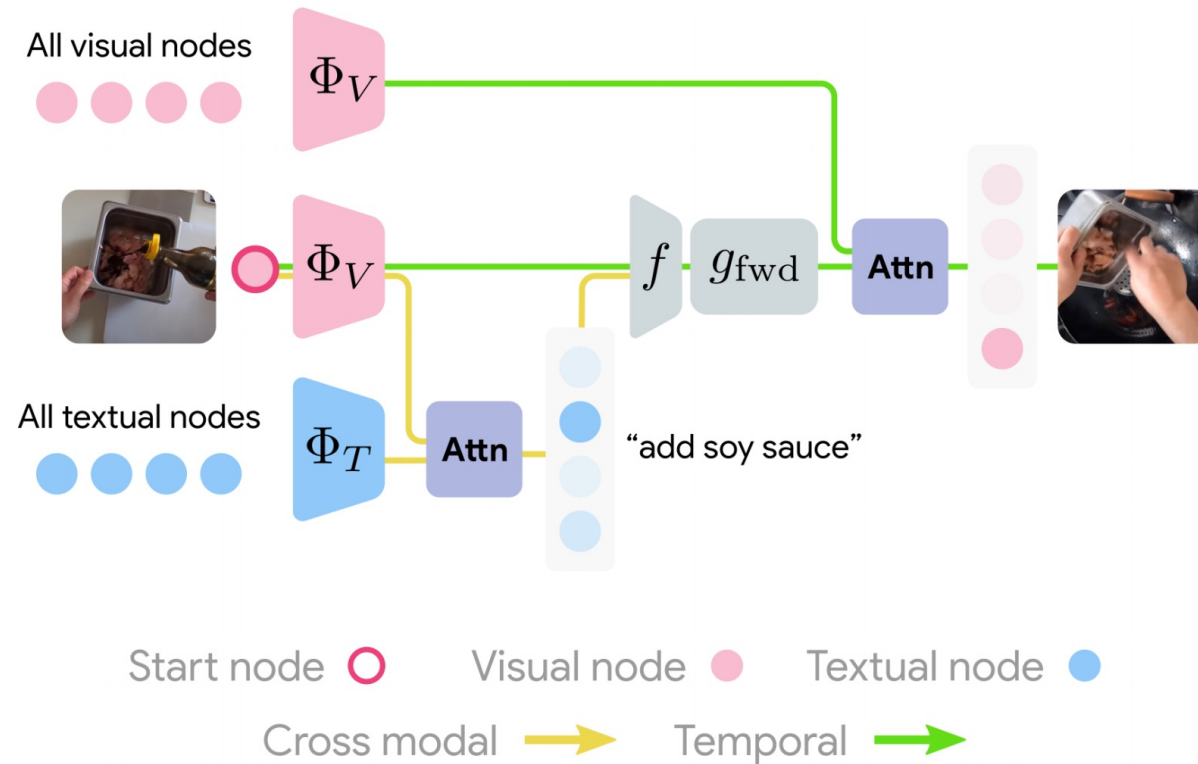(c) Time = **t+22**
"we'll be back in 30 minutes"

(d) Time = **t+35**
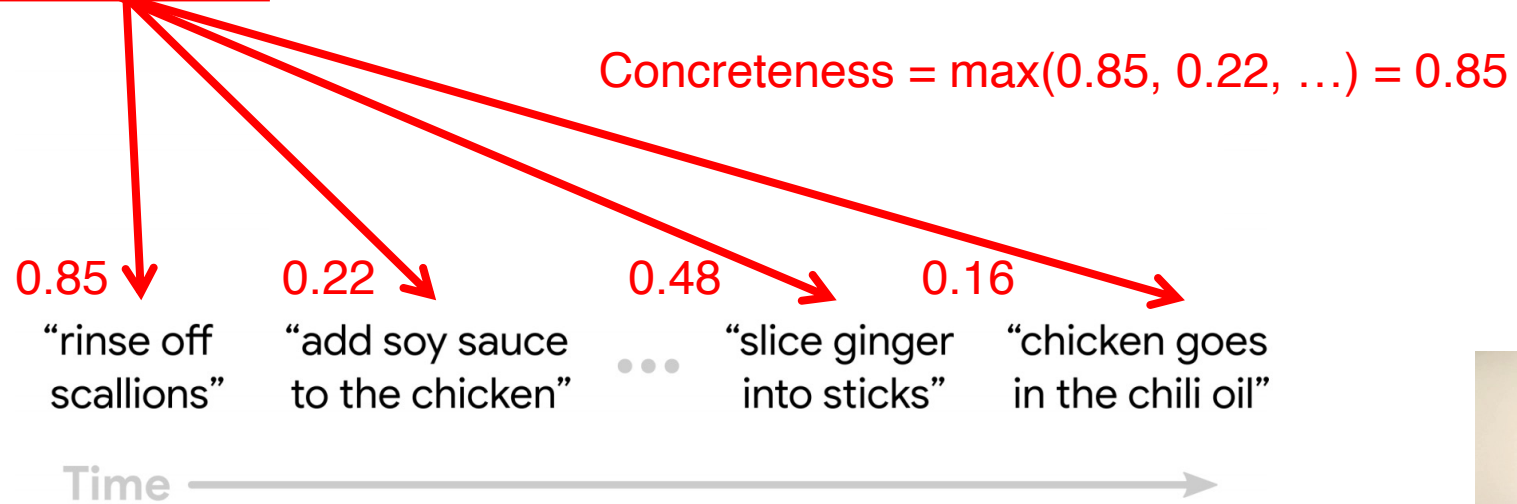"we have soft-serve ice cream"

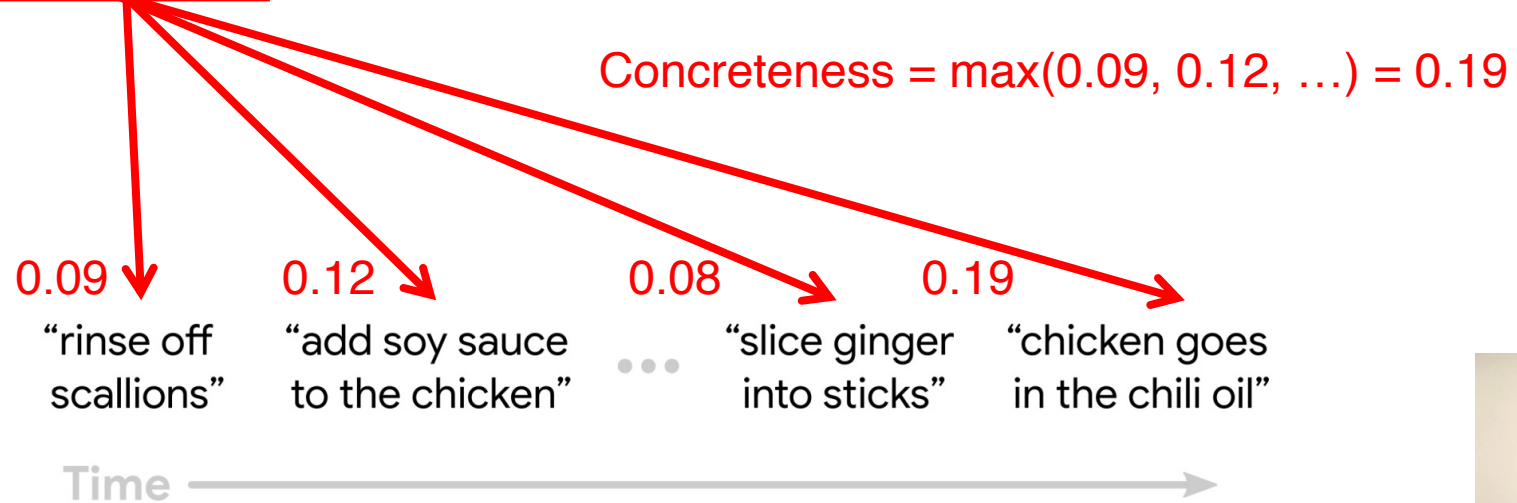# Cycling through video - implementation



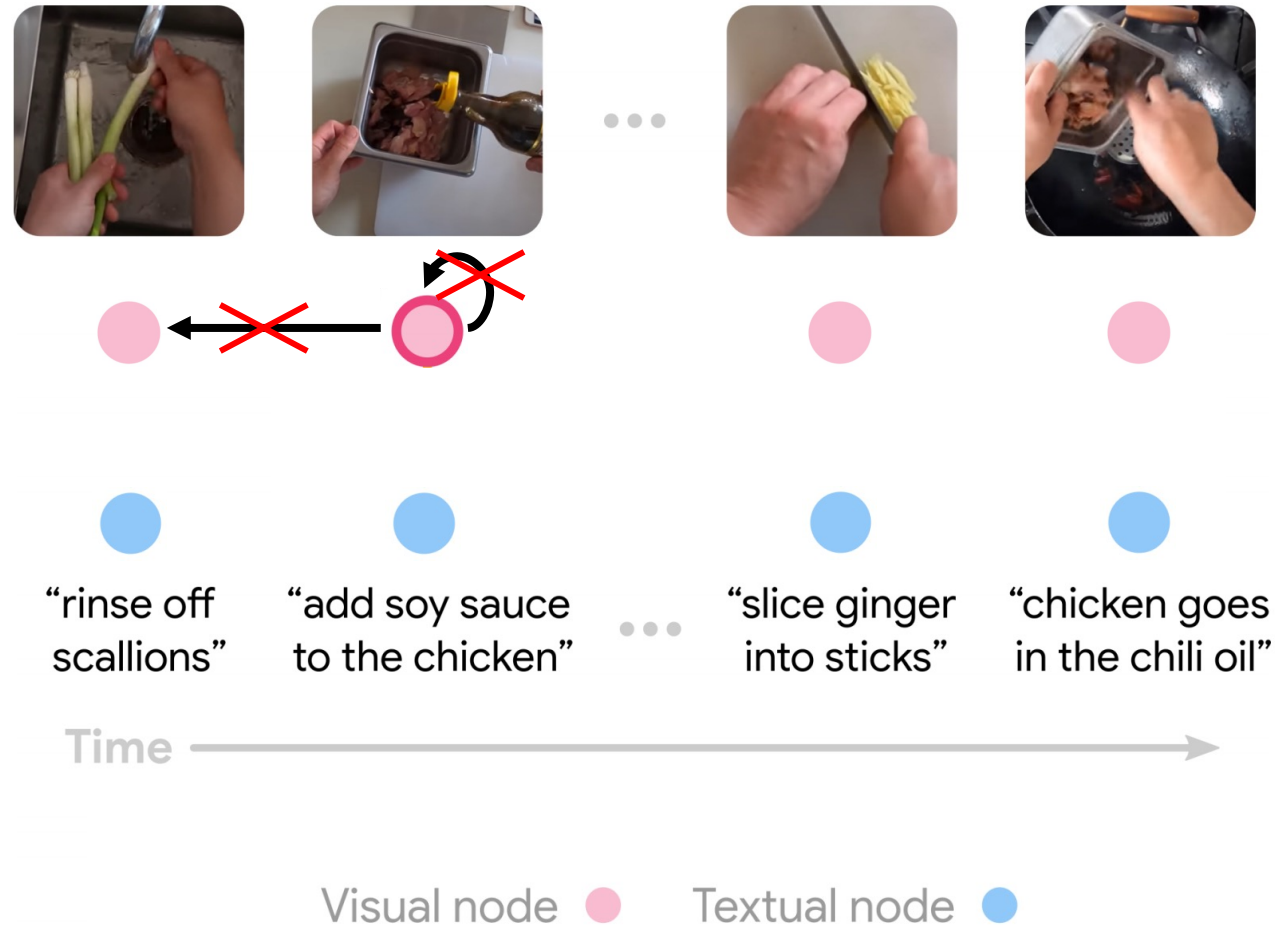Epstein, Wu, Schmid, and Sun, Learning Temporal Dynamics from Cycles in Narrated Vide

# Selecting start nodes



Concreteness = max(0.85, 0.22, …) = 0.85

0.85 "rinse off scallions"

0.22 "add soy sauce to the chicken"

0.48 "slice ginger into sticks"

0.16 "chicken goes in the chili oil"

Time

# Selecting start nodes



Concreteness = max(0.09, 0.12, …) = 0.19

0.09 "rinse off scallions"

0.12 "add soy sauce to the chicken"

0.08 "slice ginger into sticks"

0.19 "chicken goes in the chili oil"

Time

# Constraining temporal attention



"rinse off
scallions"
"add soy sauce
to the chicken"
...
"slice ginger
into sticks"
"chicken goes
in the chili oil"

Time

Visual node ● Textual node ●

# Discovering cycles in video



| Start node | Cross-modal | Forward node | Cross-modal | Backward node |
|---|---|---|---|---|
| "knead the dough until slightly sticky" | | "place dough in lightly greased bowl" | | "knead the dough until slightly sticky" |
| "get the pan hot, adding oil" | | "cook until onions are translucent" | | "get the pan hot, adding oil" |
| "pour into graham cracker crust" | | "place strawberries half inch from edge" | | "pour into graham cracker crust" |

# Finding cycles

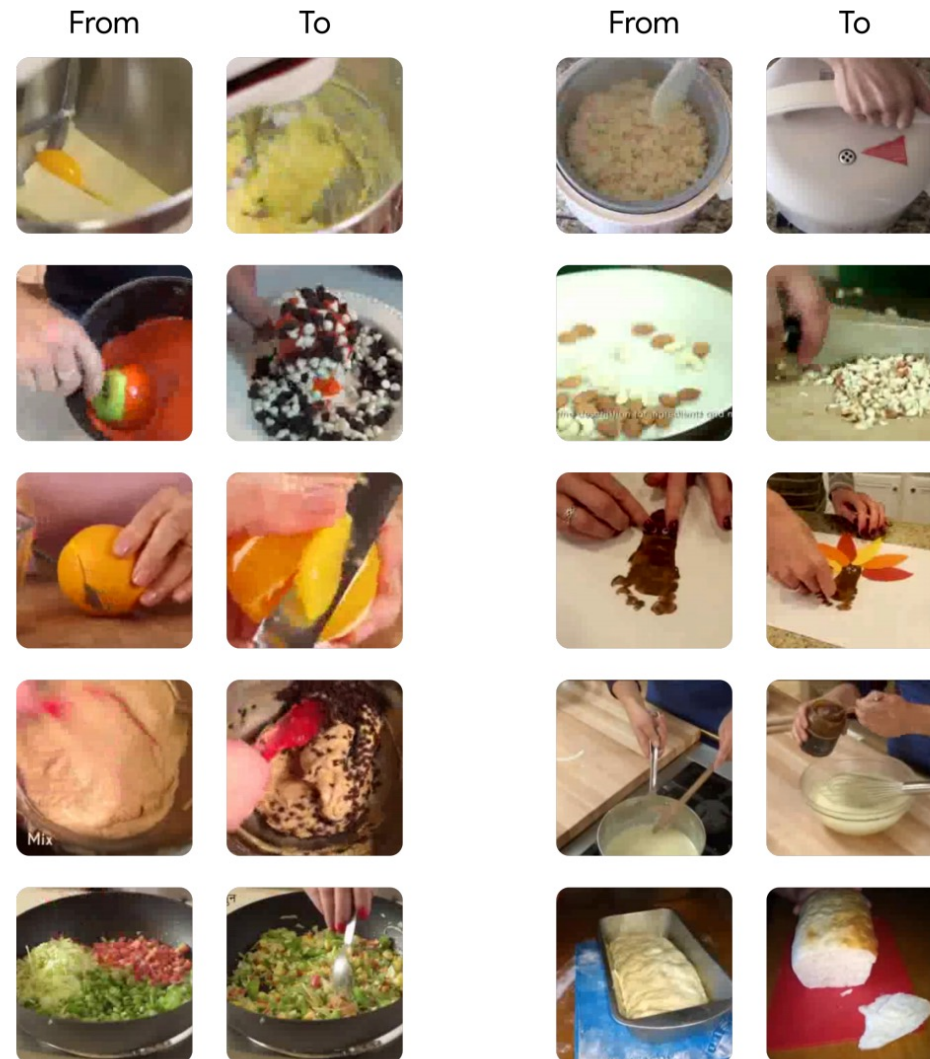| Start node | Cross-modal | Forward node | Cross-modal | Backward node |
|---|---|---|---|---|
|  | "spoon the batter into the loaf" |  | "bake until toothpick comes out clean" |  |
|  | "add the diced tomatoes" |  | "give it a quick stir to combine" |  |
|  | "cream butter in a large bowl" |  | "scoop batter into liners" |  |

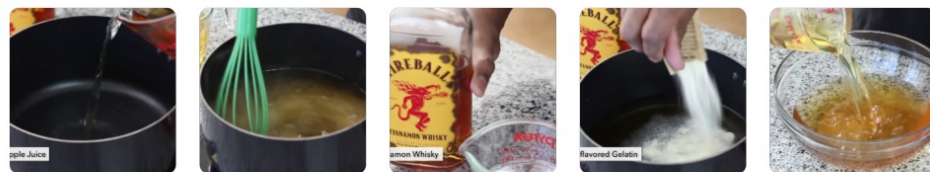# Discovering transitions in video
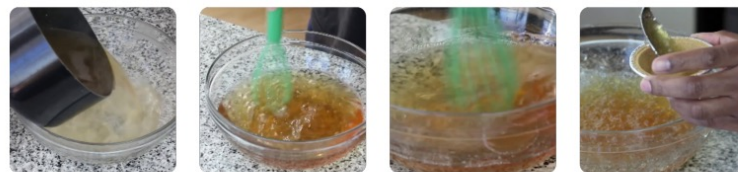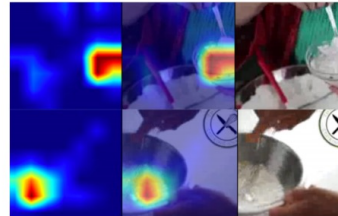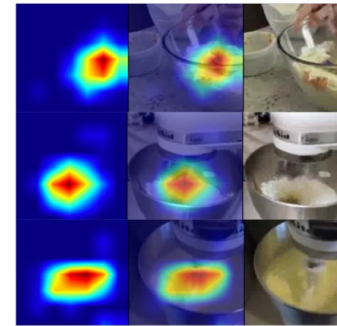
# Temporally ordering image collections

# Action and object neurons emerge
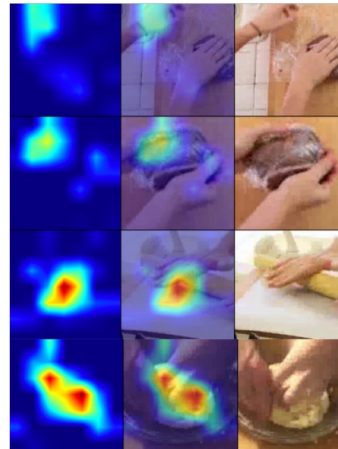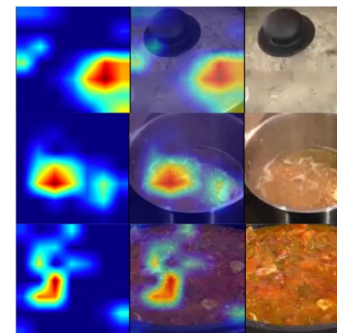


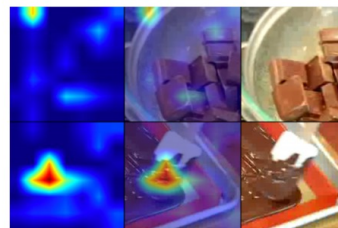flour neuron (ρ=0.172)

mix neuron (ρ=0.155)

dough neuron (ρ=0.164)

cut neuron (ρ=0.150)

chocolate neuron (ρ=0.147)

boil neuron (ρ=0.131)

# Vision-Language Navigation

ALFRED



Goal: "Rinse off a mug and place it in the coffee maker"

ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks; Shridhar et al., 2019

# VLN as a Benchmark

- Natural testbed for multimodal representations

  - Joint model visual observations, language instructions, etc.

  - From passive observation to active exploration

- The Transfer Learning Game

  - What to teach an agent before entering an environment?

  - Language and object grounding

  - Not always ideal to learn "end-to-end" and "from scr

# Focus One: language representations

$x_{1:L}$ — move to the large black end table against the wall
pick up the phone sitting on top of the end table with the blue case
carry the phone to the foot of the bed
place the phone on the bed to the right of the cushion

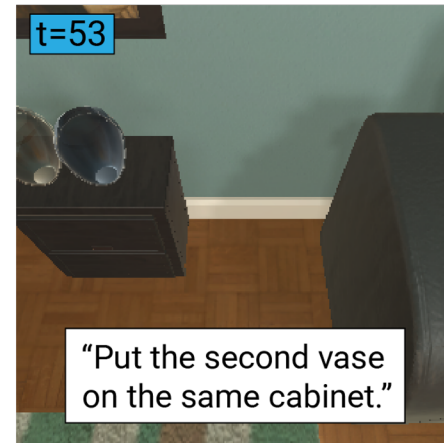$y_{1:M}$ — goto table pickup cellphone
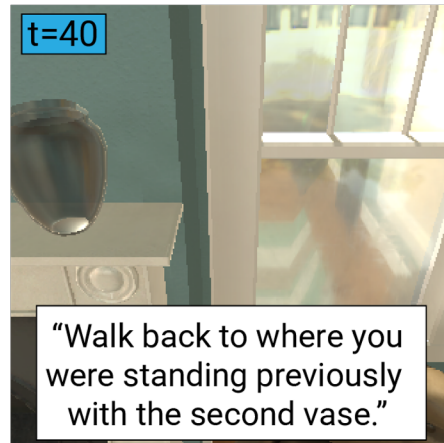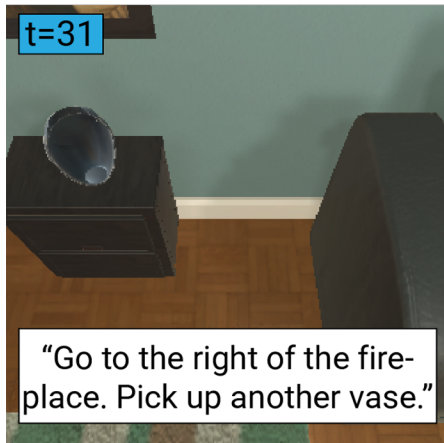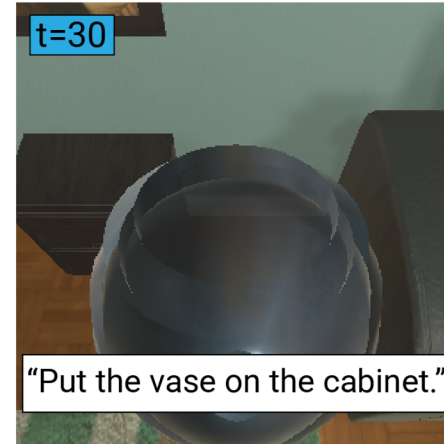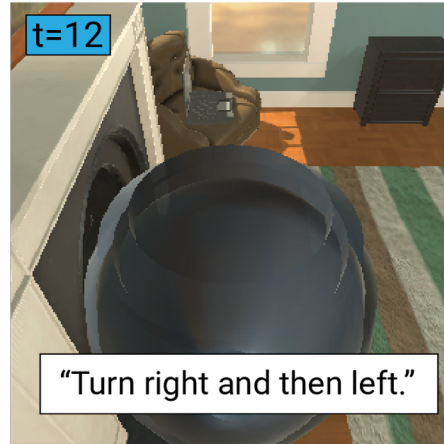goto bed put cellphone bed

Often easier to collect

Can be "pre-trained" without a specific environment.

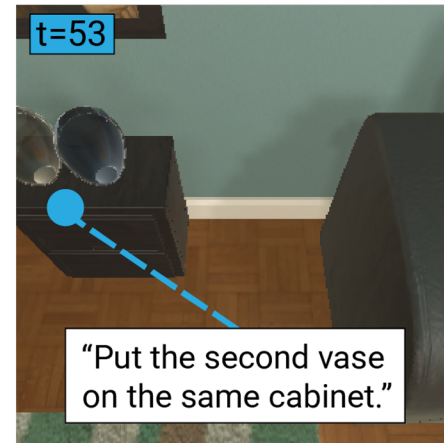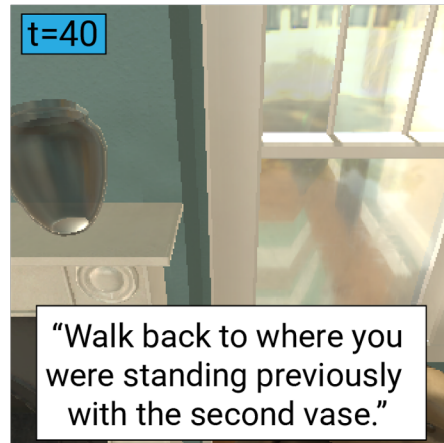Pashevich, Schmid, and Sun, Episodic Transformer for Vision-and Language Navigation,
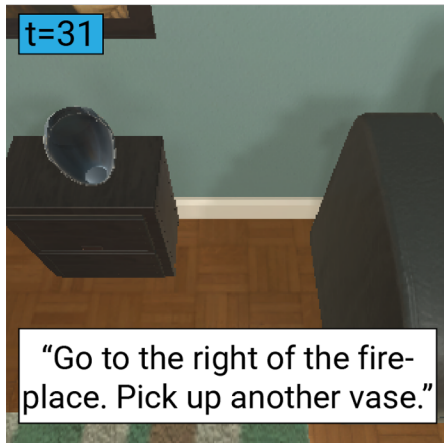
# Focus Two: Long-term dependencies

Goal: "put two vases on a cabinet"

# Focus Two: Long-term dependencies

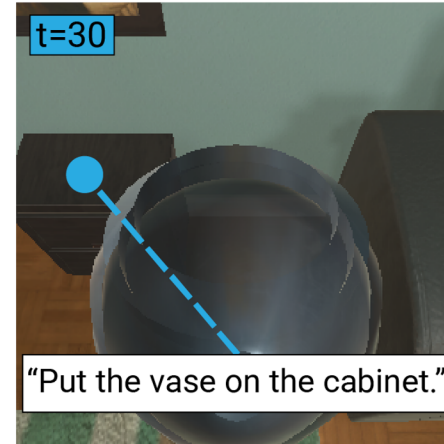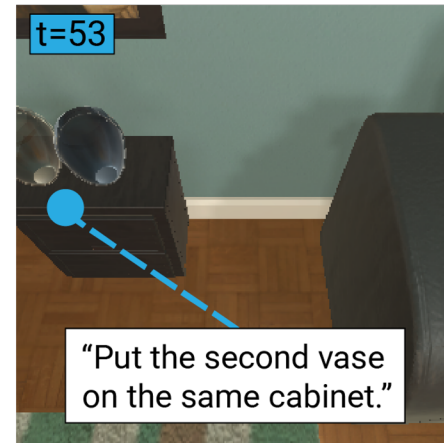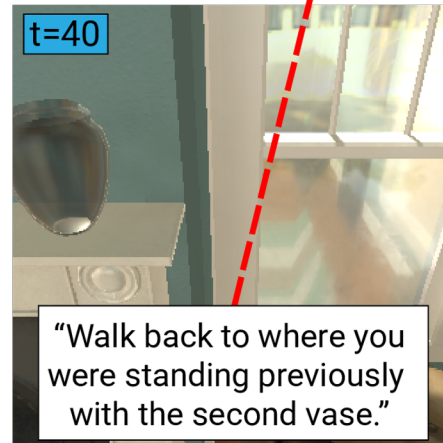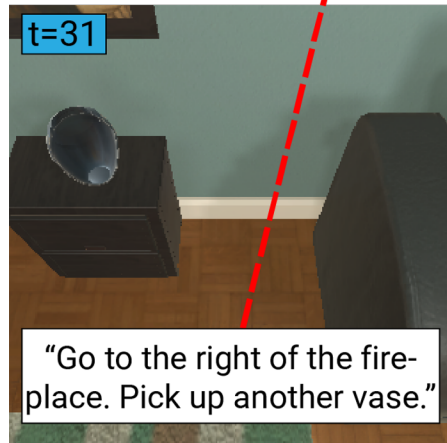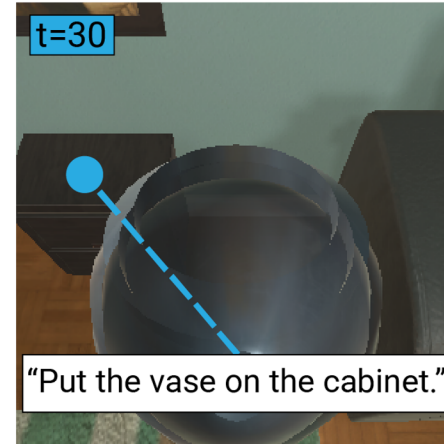Goal: "put two vases on a cabinet"

# Focus Two: Long-term dependencies

Goal: "put two vases on a cabinet"

# Results: comparison with state-of-the-art

| Model | Validation | | Test | |
|---|---|---|---|---|
| | Seen | Unseen | Seen | Unseen |
| Shridhar *et al.* [50] | 3.70 | 0.00 | 3.98 | 0.39 |
| Nguyen *et al.* [58] | N/A | N/A | 12.39 | 4.45 |
| Singh *et al.* [52] | 19.15 | 3.78 | 22.05 | 5.30 |
| E.T. (ours) | 33.78 | 3.17 | 28.77 | 5.04 |
| E.T. (ours) + synth. data | **46.59** | **7.32** | **38.42** | **8.57** |
| Human | - | - | - | 91.00 |

Comparison with state-of-the-art models.

# Self-attention to capture long-term dependency
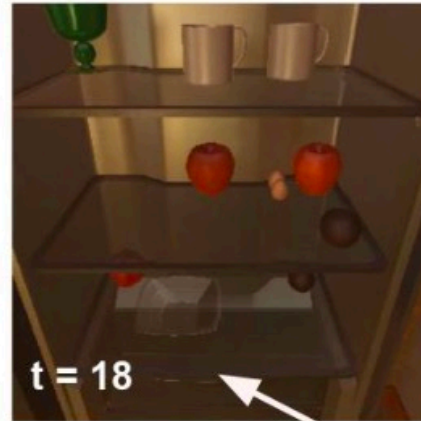


**Previous visual frames:**

t = 8
*the agent walked past a microwave*

t = 18
*the agent opened a fridge*

**Current observation:**

t = 19
*the agent needs to bring the apple back to the microwave*

**Attention to previous frames:**

**Goal:** Grab an apple, cook it and put it in the sink. **Instructions:** Turn to your left twice so that are facing the fridge. Open the fridge, grab an apple from the shelf and close the fridge door. *to the left of the fridge to face the microwave.* Put the apple in the microwave and cook it for a seconds before taking it back out and closing the microwave. Turn to face your left. Put the ap in the sink.

# Code and checkpoints are released!

https://github.com/alexpashevich/E.T.

# Summary

- Many interesting tasks for detailed video understanding

  - Video is encyclopedia of multimedia contents!

- From manual annotation to "automatic" supervision

  - Self-supervised: Contrastive Learning

  - Cross-modal supervised: Cross-modal cycle consistency

- Many interesting applications of detailed video understanding

  - Structured multimodal representations for navigation

  - Better interpretable, more generalizable models

# Collaborators