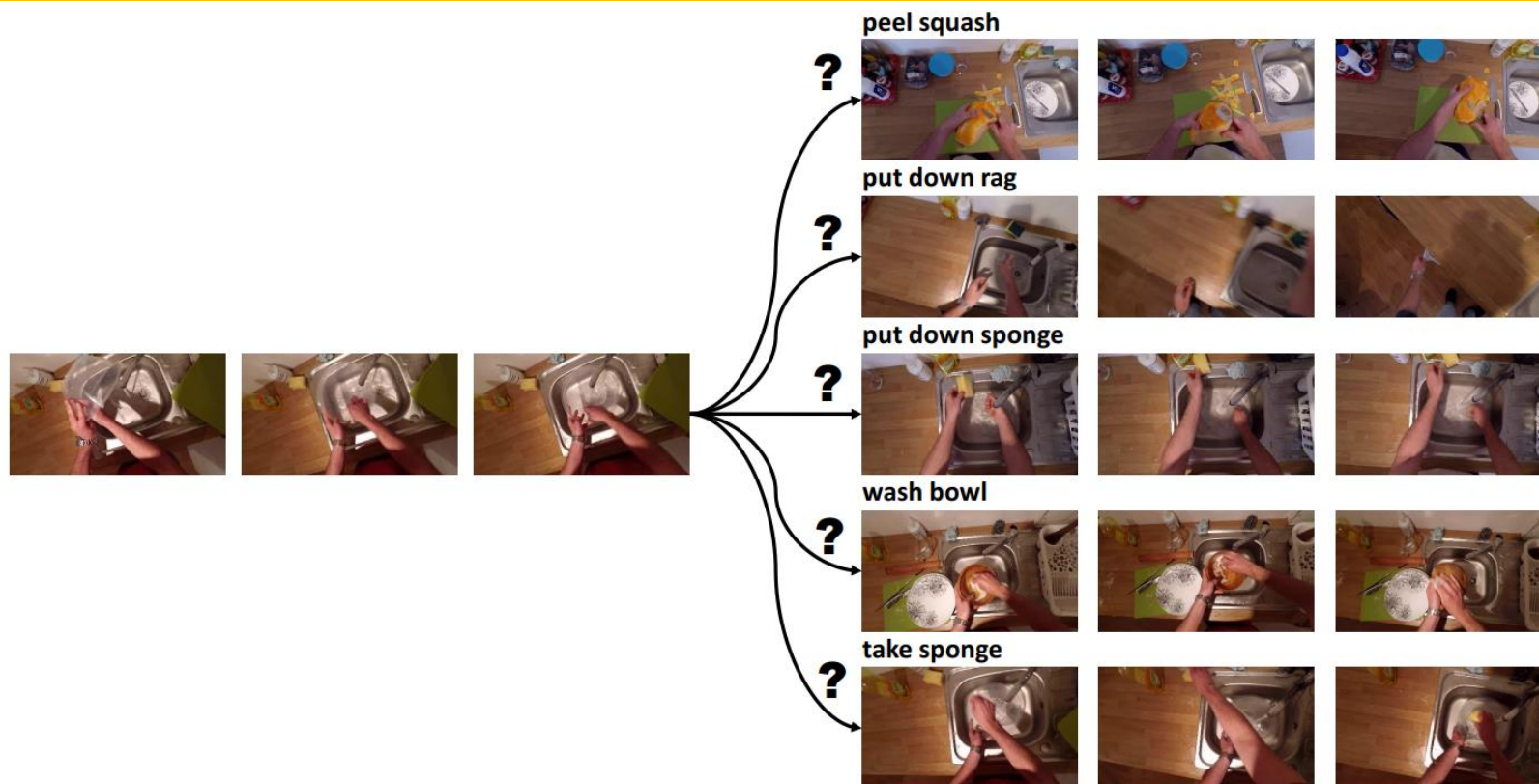# Part5: Cooking action anticipation

# Motivation

**Predicting what will happen in the future!**

# Motivation

● **Egocentric Vision**

☐ wearable cameras          ☐ daily activities



insert coffee holder

stir chicken

put tongs in sink

move tap

# Motivation

- ➢ in the kitchen environment

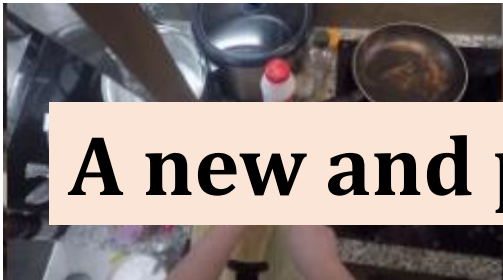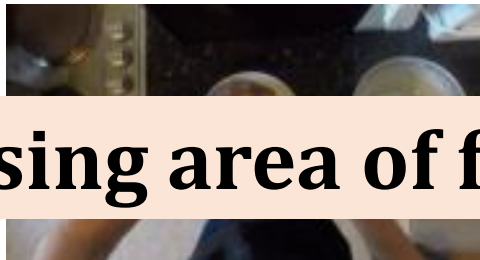- ➢ cooking-related actions



fry egg



pour milk
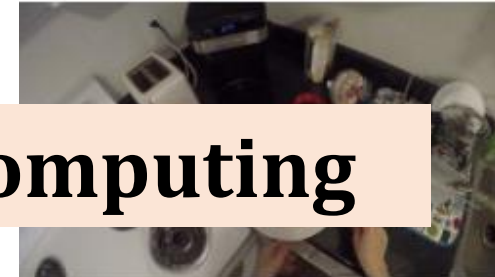


split salmon



slice chilli



flip fish



apply spreads

**A new and promising area of food computing**

# Motivation

➢ **In the kitchen environment**





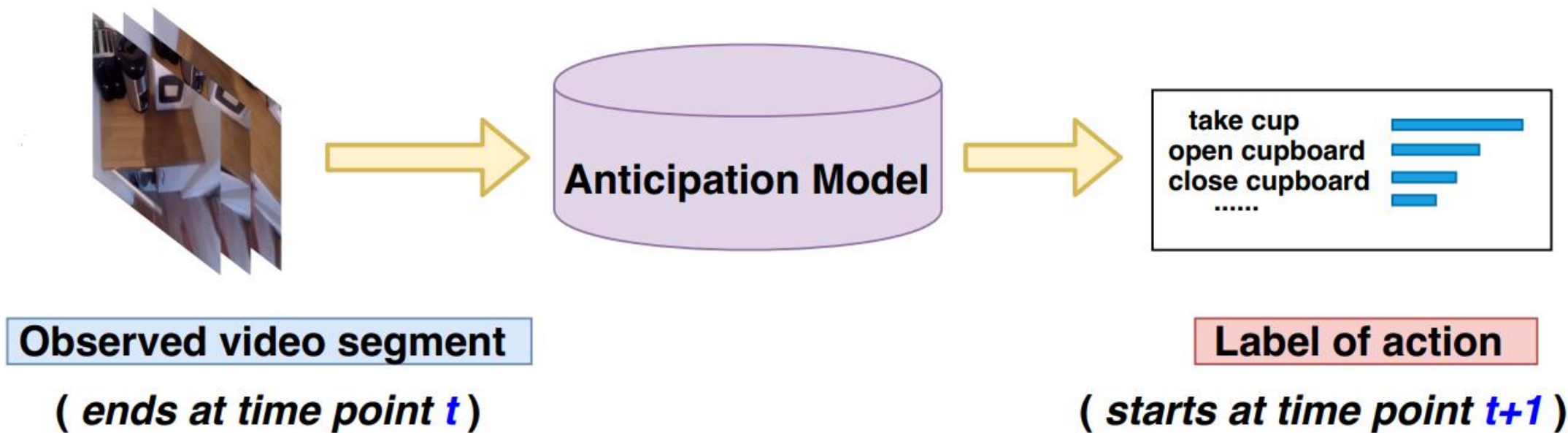➢ **Application of sevice robots**

☐ Help those who are disabled to cook recipes, wish dishes, etc
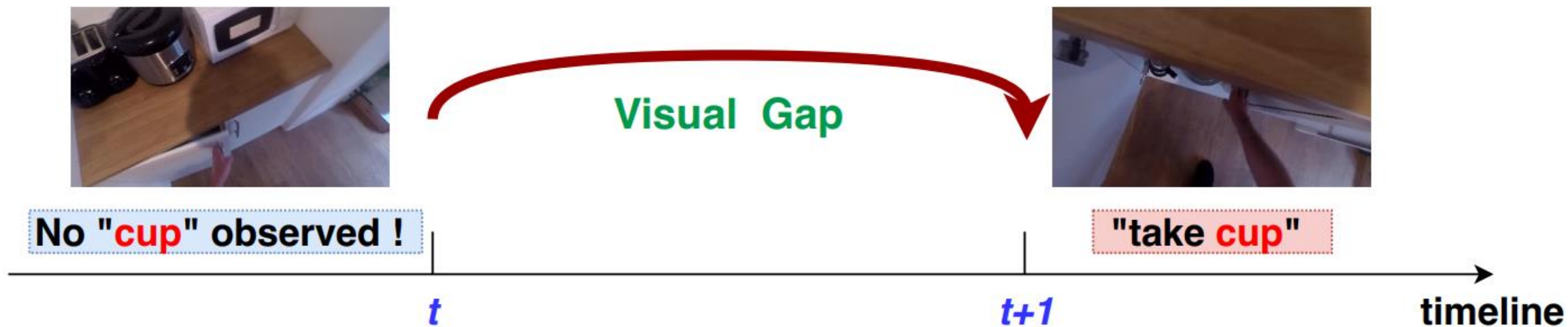
☐ Instruct people to learn how to cook
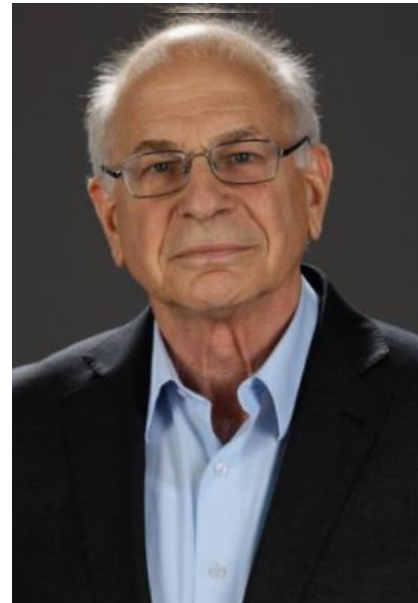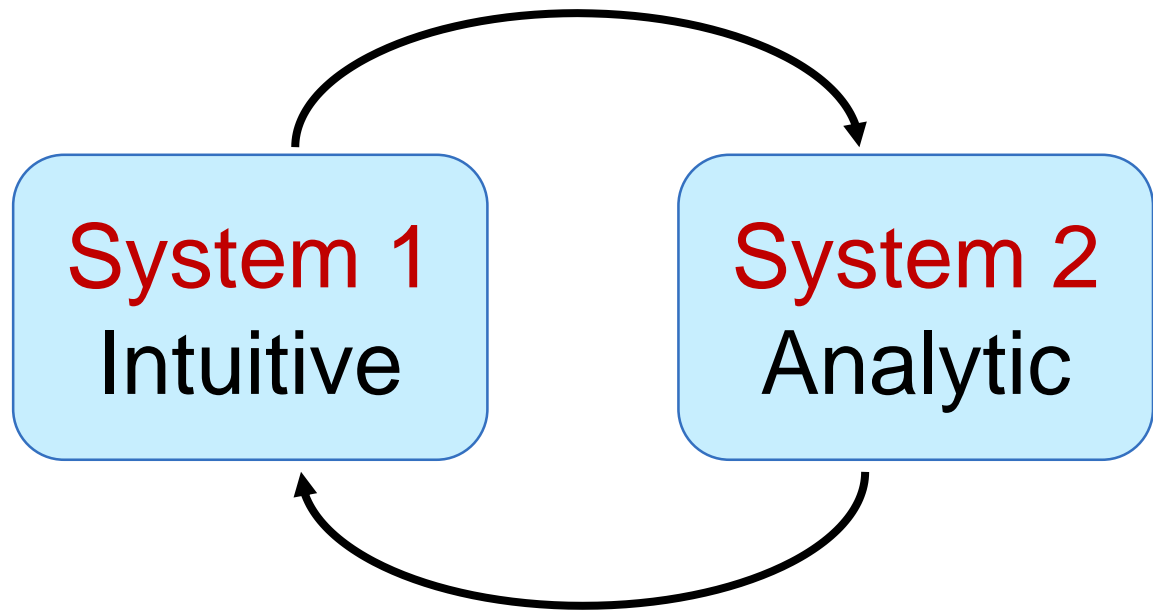
# Egocentric Action Anticipation

➢ **Definition**
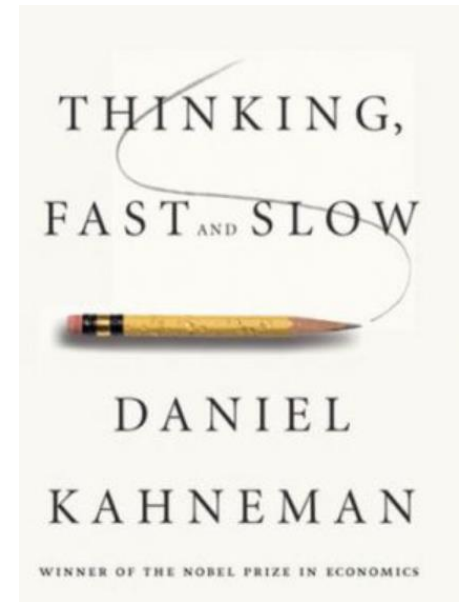
# Egocentric Action Anticipation

➢ **Difficulty**

# Exploration from human psychology

☐ Two modes in the cognitive system of human brain: **intuition and analysis**

☐ Intuition and anlysis are both crucial in solving many problems (e.g. making predictions)
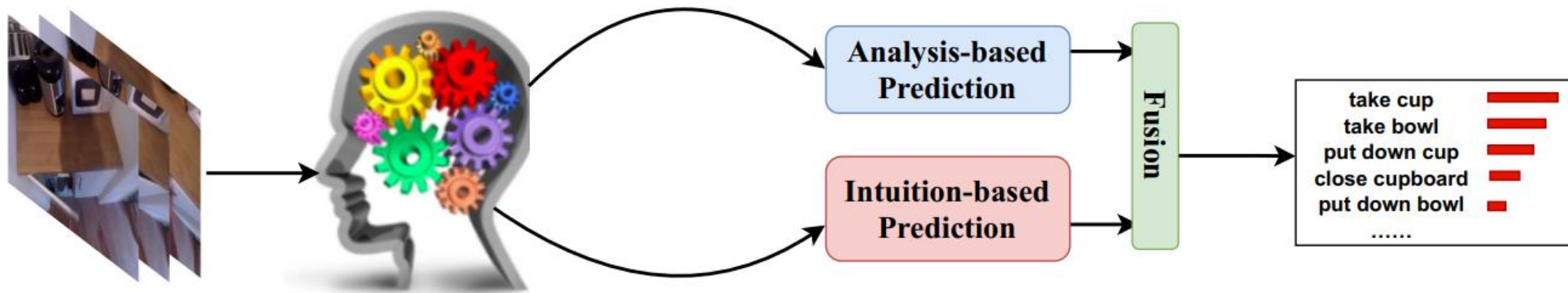


Daniel Kahneman

8

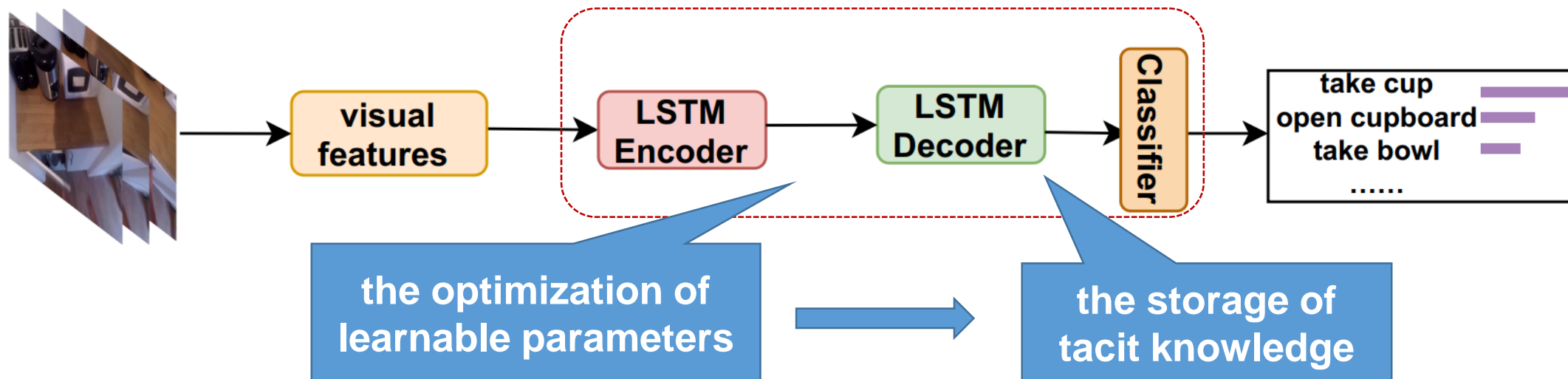# Exploration from human psychology

Construct a basic framework that integrates both intuition-based prediction and analysis-based prediction to imitate human beings in making predictions

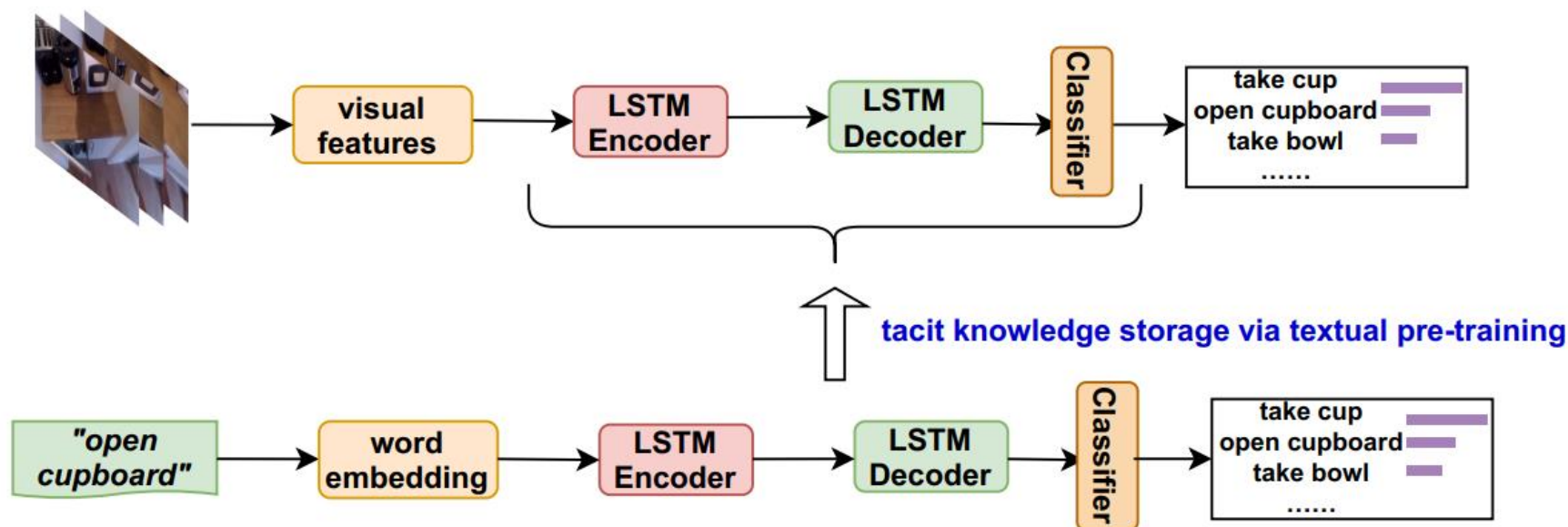# Exploration from human psychology

☐ Intuition-based prediction

➢ Subconscious, habitual

➢ Tacit knowledge (**hard to explain**)

➢ An **encoder-decoder** structure (a black-box process)



10

# Exploration from human psychology

❑ Intuition-based prediction

➤ Visual information is insufficient to store tacit knowledge

➤ Introduce textual pre-training to store tacit knowledge in advance

# Exploration from human psychology

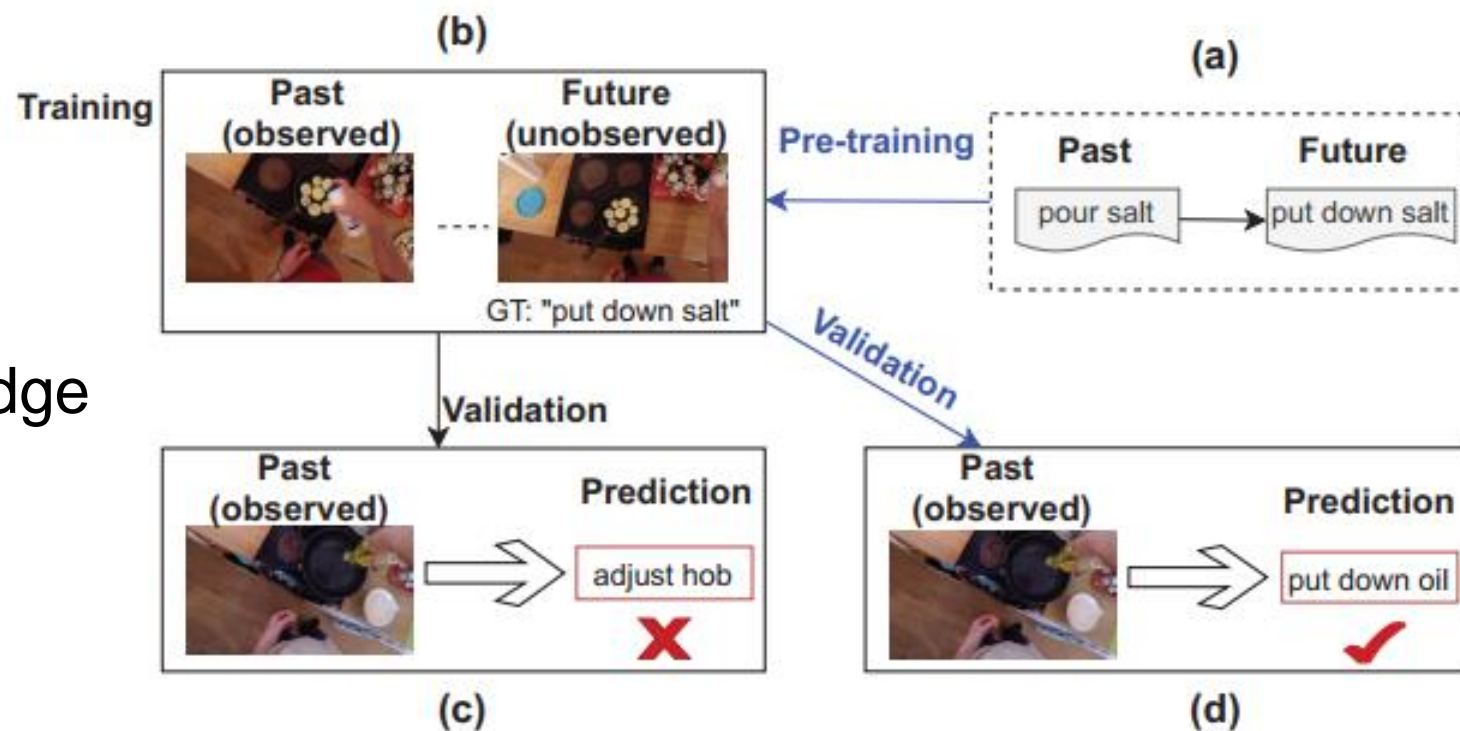☐ Intuition-based prediction

➢ (b)→(c):

☐ Visual information

☐ Insufficient to store tacit knowledge
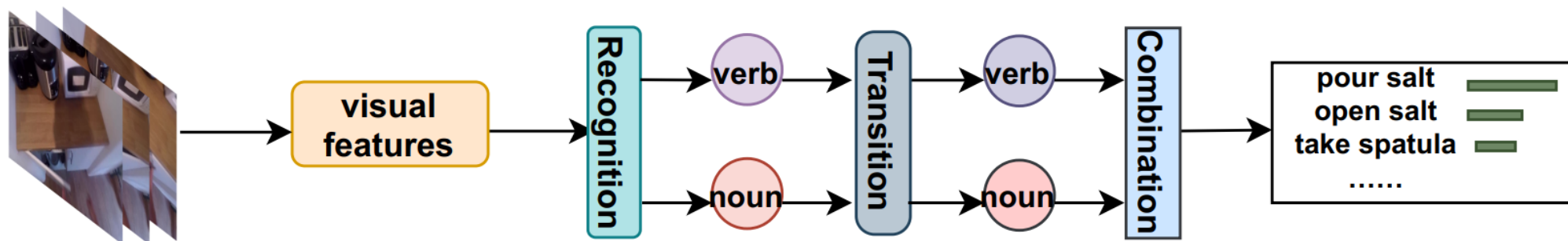
➢ (a)→(b)→(d):

☐ Visual + text information

☐ Store more reliable tacit knowledge

# Exploration from human psychology
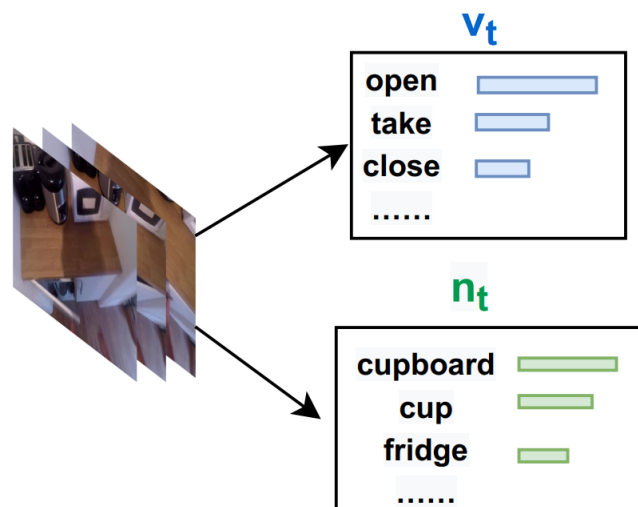
## ❑ Analysis-based prediction

➢ Conscious and explicit

➢ Tend to process information under given principles
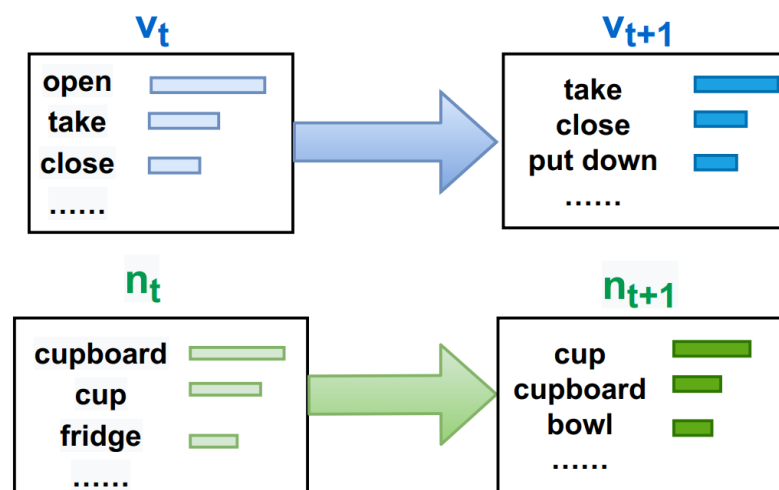
➢ An interpretable three-step pipeline

# Exploration from human psychology
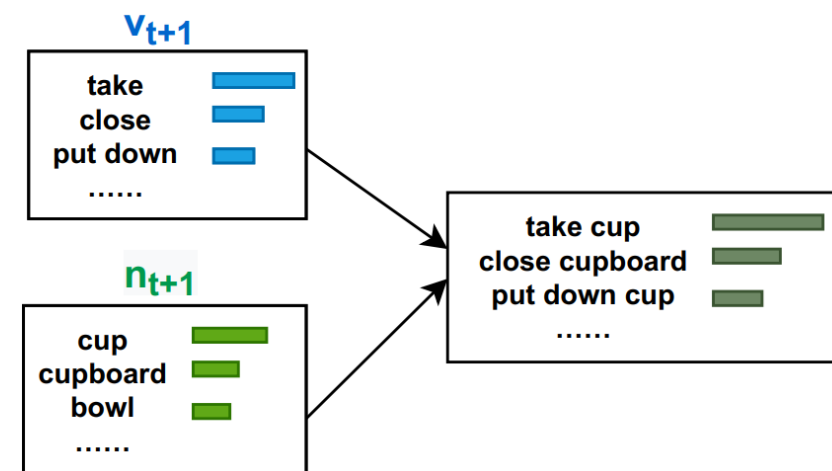
☐ Analysis-based prediction
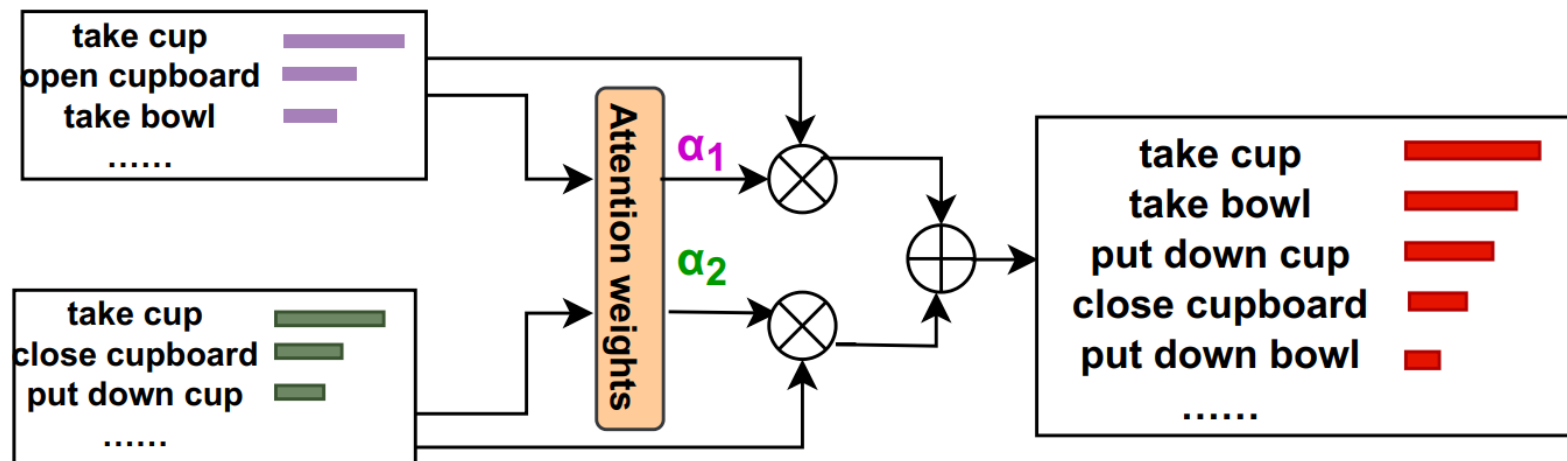


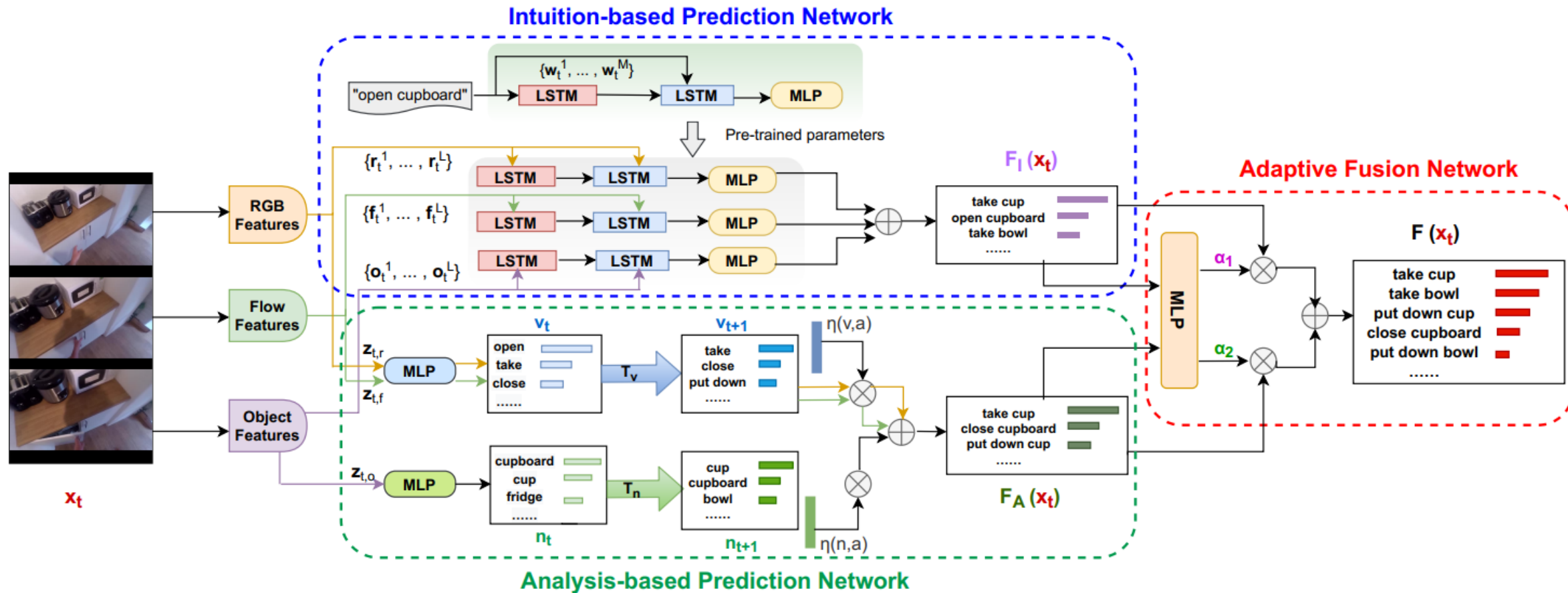➢ **Recognition**   ➢ **Transition**   ➢ **Combination**

# Exploration from human psychology

☐ Intuition-analysis fusion

➢Both intuition and analysis are crucial and indispensable

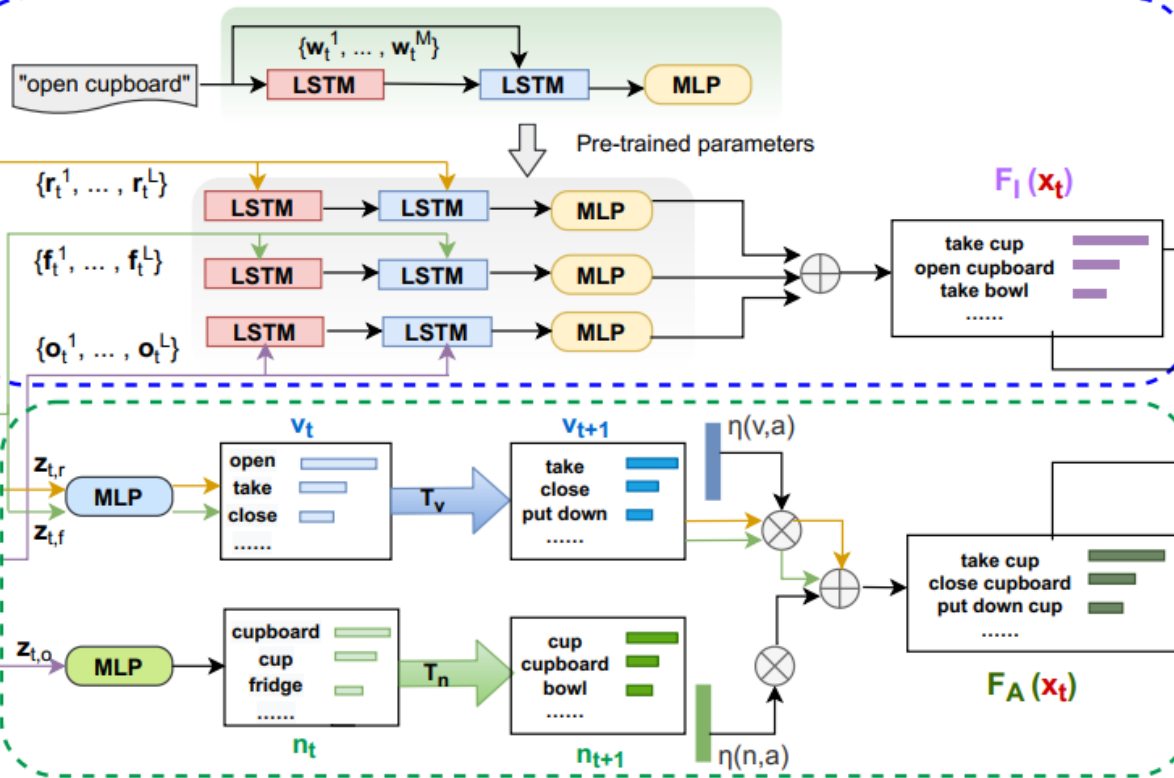➢Compute attention weights for intuition-based and analysis-based prediction and integrate them adaptively

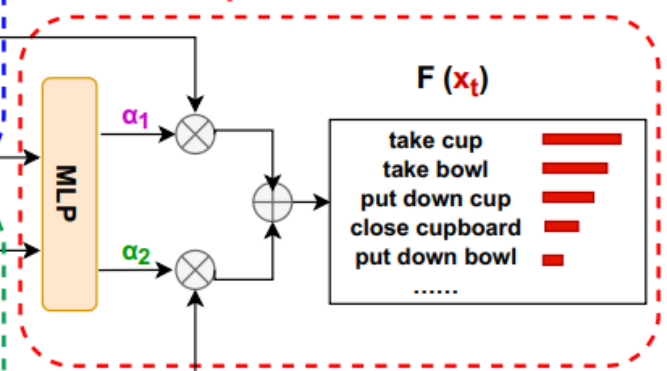# Intuition-Analysis Integrated Framework

# Intuition-Analysis Integrated Framework



LSTM:
store tacit knowledge

Intuition-based Prediction Network

$\{w_t^1, \ldots, w_t^M\}$

"open cupboard" → LSTM → LSTM → MLP

Pre-trained parameters

$\{r_t^1, \ldots, r_t^L\}$ LSTM → LSTM → MLP

$\{f_t^1, \ldots, f_t^L\}$ LSTM → LSTM → MLP

$\{o_t^1, \ldots, o_t^L\}$ LSTM → LSTM → MLP

$F_I(x_t)$

take cup
open cupboard
take bowl
......

RGB Features

Flow Features

Object Features

$x_t$

$z_{t,r}$ MLP
$z_{t,f}$

$v_t$
open
take
close
......

$T_v$

$v_{t+1}$
take
close
put down
......

$\eta(v,a)$

$z_{t,o}$ MLP

$n_t$
cupboard
cup
fridge
......

$T_n$

$n_{t+1}$
cup
cupboard
bowl
......

$\eta(n,a)$

take cup
close cupboard
put down cup
......

$F_A(x_t)$

Analysis-based Prediction Network

Markov logic:
transit from past to future

Adaptive Fusion Network

MLP
$\alpha_1$
$\alpha_2$

$F(x_t)$

take cup
take bowl
put down cup
close cupboard
put down bowl
......

Fusion:
$$F(x_t) = a_1 * F_I(x_t) + a_2 * F_A(x_t)$$

# Evaluation

## ☐ **EPIC-Kitchens Dataset**

- ➤ 32 kitchens - 4 cities

- ➤ Head-mounted camera

- ➤ 55 hours of recording - Full HD, 60fps

- ➤ 11.5M frames

- ➤ 39,594 action segments

- ➤ 125 verb classes

- ➤ 352 noun classes

- ➤ 2,513 action classes

# Evaluation

## ☐ Comparison with other methods

| Method | Top1@V | Top1@N | Top1@A | Top5@V | Top5@N | Top5@A |
|---|---|---|---|---|---|---|
| 2SCNN [30] | 25.23 | 9.97 | 2.29 | 68.66 | 27.38 | 9.35 |
| TSN [35] | 25.30 | 10.41 | 2.39 | 68.32 | 29.50 | 9.63 |
| TSN+MCE [9] | 21.27 | 9.90 | 5.57 | 63.33 | 25.50 | 15.71 |
| Miech *et al.* [20] | **28.37** | 12.43 | 7.24 | 69.96 | 32.20 | 19.29 |
| RULSTM [10] | 27.01 | **15.19** | 8.16 | 69.55 | 34.38 | 21.20 |
| Ours-IPN | 27.24 | 14.58 | 8.06 | 69.17 | 34.21 | 20.21 |
| Ours-APN | 24.07 | 14.65 | 7.27 | 68.62 | 34.45 | 18.33 |
| Ours-IAI | 27.89 | 14.89 | **8.57** | **70.06** | **35.51** | **21.41** |

➤ Top1@V/N/A: Top-1 accuracy for verbs/nouns/actions

➤ Top5@V/N/A: Top-5 accuracy for verbs/nouns/actions

# Evaluation



Sucessful cases: intuition and analysis complement each other

More food-related information from cooking domain is needed !
(e.g.,  ingredient information of the being-prepared dishes)