# Temporal AND-OR Graph
## for representation and recognition of Events, Actions, Motions

Song-Chun Zhu, Sinisa Todorovic, and Ales Leonardis

At CVPR, Providence, Rhode Island
June 16, 2012

1

# Goal: Recognize events in daily scenes

▸ For example, an office.

Ref.   Pei, Si and Zhu,  ICCV2011.

# Challenges

1. Events happen over an extended time period

   - Variant time-span

   - Could be interrupted

   - Multiple routes

   - Intention and prediction



2. Actions are hard to recognize

   - Subtle and similar

   - No salient motion/pose at most of the time

   - Contextual objects -- key!!
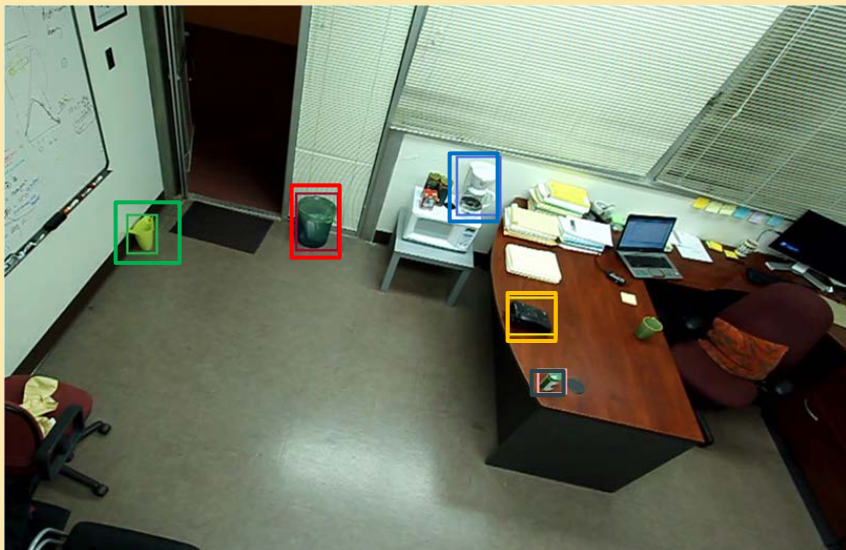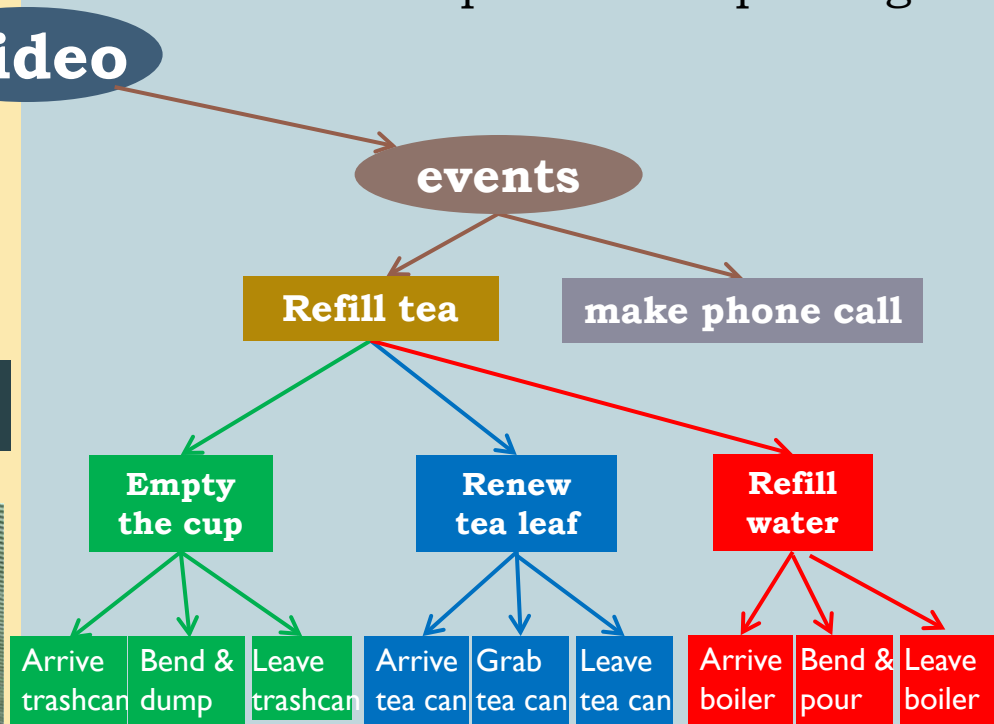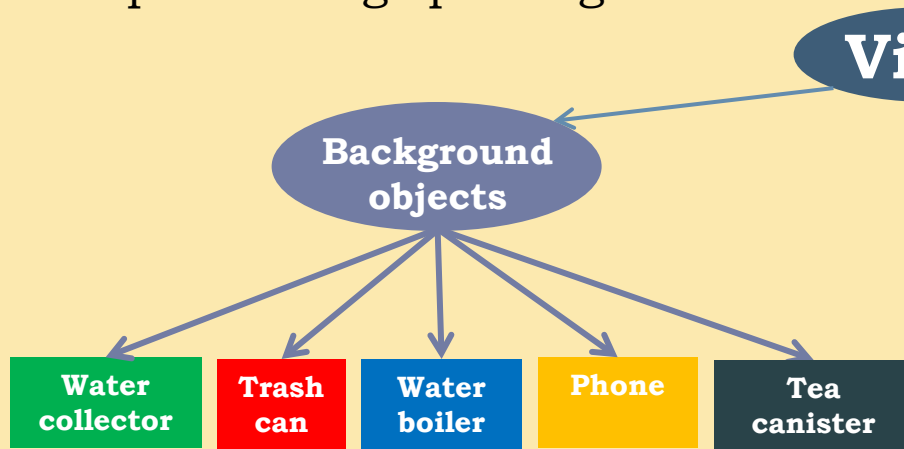


Use laptop    Read book    Dump water    Use microwave

# Overview of our approach

Spatial image parsing

Temporal event parsing

# Scene parsing

# How to define actions and events?

Some neurons in the pre-motor area encode actions



**Mirror neurons** firing when performing action or seeing other people performing the action

Gallese et al 96, Rizzolatti et a. 96

# Actions = Spatiotemporal relations between body parts and objects in the scene

| Atomic Actions | Fluents | Symbols | | Examples |
|---|---|---|---|---|
| | | Foreground | Background | |
| Shake Hands(P1,P2) | Near(P1,P2) And Touch (P1.hand, P2.hand) | | | |
| Use Dispenser(P3) | Bend(P3) and Near(P3, A) And Touch(P3.hand, A) | | | |
| Pick up Phone(P4) | Touch(P4, B) And On(B) | | | |

Some of the learned atomic actions by pursuing the co-occurrence of relations.

M.T. Pei, Z.Z. Si, B. Yao, and S.C. Zhu, "Video Event Parsing and Learning with Goal and Intent Prediction," 2012

# Actions = Spatiotemporal relations between body parts and objects in the scene

## Unary Relations

| Status of person | Symbols | Examples | Status of objects | Examples |
|---|---|---|---|---|
| Stand(P1) | | | On (phone) | |
| Stretch(P1) | | | Off (phone) | |
| Bend (P1) | | | On(screen) | |
| Sit (P2) | | | Off(screen) | |

## Binary Relations

| Binary Fluent (A,B) | Touch (A,B) | Near (A,B) | Occlude (A,B) | In(A,B) |
|---|---|---|---|---|
| Symbols | A B | A B | A B | A B |
| Examples | | | | |

High-order relations:
- E.g., surrounded by

8

# Event as temporal And-Or-Graph

# Event as temporal And-Or-Graph



And-node

Or-node

Leaf-node

Make phone call

Get close to phone

Pick up receiver

Talking over phone

Put down receiver

Stand still

Walk around

Get close to phone

Pick up phone

Stand still

Walk around

Put down receiver

# Formulation



$$p(g) = \frac{1}{Z}exp\{score(g)\}$$

Grammar

| Data term | Or node Frequency term | Temporal Relations |
|---|---|---|

$$score(g) = \sum_{v_t \in T(g)} \lambda_{v_t}\alpha(v_t) + \sum_{v \in V_o(g)} \lambda_v\omega(v) + \sum_{(i,j) \in E(g)} \lambda_{ij}r_{ij}(v_i, v_j)$$

# Formulation



$$p(g) = \frac{1}{Z} exp\{score(g)\}$$

Grammar

Data term | Or node Frequency | Temporal Relations

$$score(g) = \sum_{v_t \in T(g)} \lambda_{v_t} \alpha(v_t) + \sum_{v \in V_o(g)} \lambda_v \omega(v) + \sum_{(i,j) \in E(g)} \lambda_{ij} r_{ij}(v_i, v_j)$$

$$\alpha(v_t) = \sum_{i \in \mathcal{F}} \beta_i, h_i(v_t) - dist(P_{person}, P_{obj})$$

12

# Parsing process (Earley Parser [Earley 1970])

# Parsing process (Earley Parser [Earley 1970])

# Parsing process (Earley Parser [Earley 1970])

# Parsing process (Earley Parser [Earley 1970])



And-node

Or-node

Leaf-node

Refill tea

Empty cup

Renew tea leaf

Refill water

Dump old tea

Dump water

Null

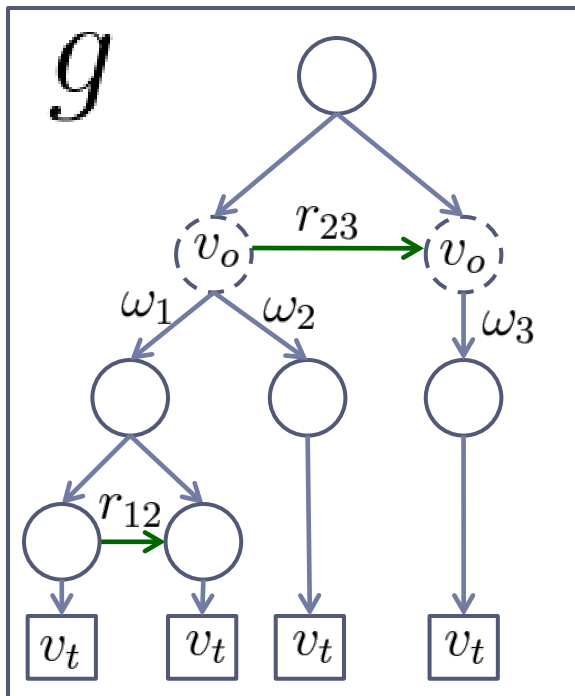Renew tea leaf

Refill water

Arrive trashcan

Bend down

Leave trashcan

Arrive collector

Bend down

Leave collector

NULL

Get close to tea can

Grab tea can

Leave tea-can

Arrive boiler

Bend down

Leave boiler

16

# Parsing process (Earley Parser [Earley 1970])

# Parsing process (Earley Parser [Earley 1970])

# Parsing process (Earley Parser [Earley 1970])

# Parsing process (Earley Parser [Earley 1970])

# Parsing process (Earley Parser [Earley 1970])



And-node
Or-node
Leaf-node

Refill tea

Empty cup          Renew tea leaf          Refill water

Dump old tea       Dump water       Null       Renew tea leaf       Refill water

Arrive trashcan | Bend down | Leave trashcan | Arrive collector | Bend down | Leave collector | NULL | Get close to tea can | Grab tea can | Leave tea-can | Arrive boiler | Bend down | Leave boiler

# Parsing process (Earley Parser [Earley 1970])



And-node
Or-node
Leaf-node

Refill tea

Empty cup    Renew tea leaf    Refill water

Dump old tea    Dump water    Null    Renew tea leaf    Refill water

Arrive trashcan    Bend down    Leave trashcan    Arrive collector    Bend down    Leave collector    NULL    Get close to tea can    Grab tea can    Leave tea-can    Arrive boiler    Bend down    Leave boiler

# Parsing process (Earley Parser [Earley 1970])

# Parsing process (Earley Parser [Earley 1970])



And-node

Or-node

Leaf-node

Refill tea

Empty cup

Renew tea leaf

Refill water

Dump old tea

Dump water

Null

Renew tea leaf

Refill water

Arrive trashcan

Bend down

Leave trashcan

Arrive collector

Bend down

Leave collector

NULL

Get close to tea can

Grab tea can

Leave tea-can

Arrive boiler

Bend down

Leave boiler

# Parsing process (Earley Parser [Earley 1970])



And-node
Or-node
Leaf-node

Refill tea

Empty cup · Renew tea leaf · Refill water

Dump old tea · Dump water · Null · Renew tea leaf · Refill water

Arrive trashcan | Bend down | Leave trashcan | Arrive collector | Bend down | Leave collector | NULL | Get close to tea can | Grab tea can | Leave tea-can | Arrive boiler | Bend down | Leave boiler

# Parsing process (Earley Parser [Earley 1970])



And-node
Or-node
Leaf-node

Refill tea

Empty cup

Renew tea leaf

Refill water

Dump old tea

Dump water

Null

Renew tea leaf

Refill water

| Arrive trashcan | Bend down | Leave trashcan | Arrive collector | Bend down | Leave collector | NULL | Get close to tea can | Grab tea can | Leave tea-can | Arrive boiler | Bend down | Leave boiler |

26

# Parsing process (Earley Parser [Earley 1970])



And-node
Or-node
Leaf-node

Refill tea

Empty cup · Renew tea leaf · Refill water

Dump old tea · Dump water · Null · Renew tea leaf · Refill water

Arrive trashcan · Bend down · Leave trashcan · Arrive collector · Bend down · Leave collector · NULL · Get close to tea can · Grab tea can · Leave tea-can · Arrive boiler · Bend down · Leave boiler
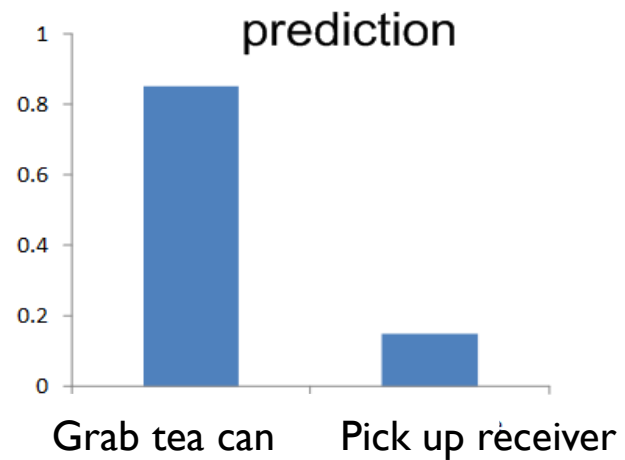
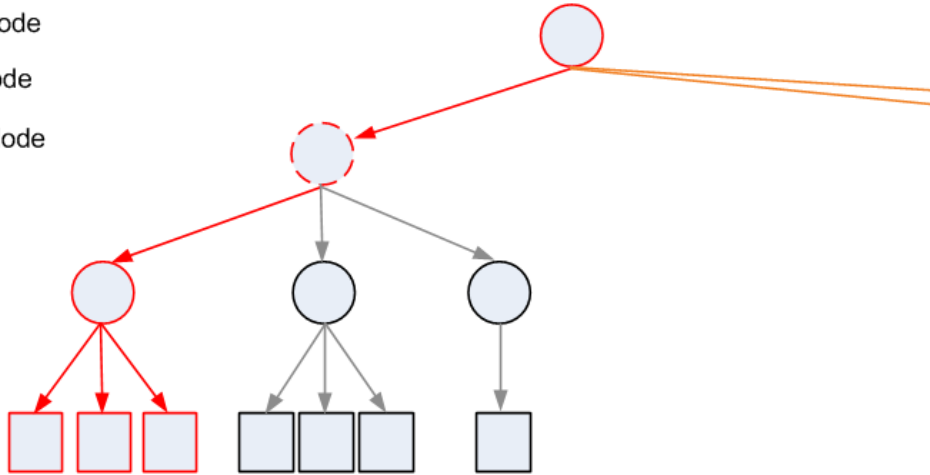# Parsing: A modified Earley parser [Earley 1970]

# Intention and prediction

# Intention and prediction



Time 2  Refill tea  Make phone call

intention
prediction
Observed action
Predicted action

intention
Refill tea  Make phone call

prediction
Grab tea can  Pick up receiver

30

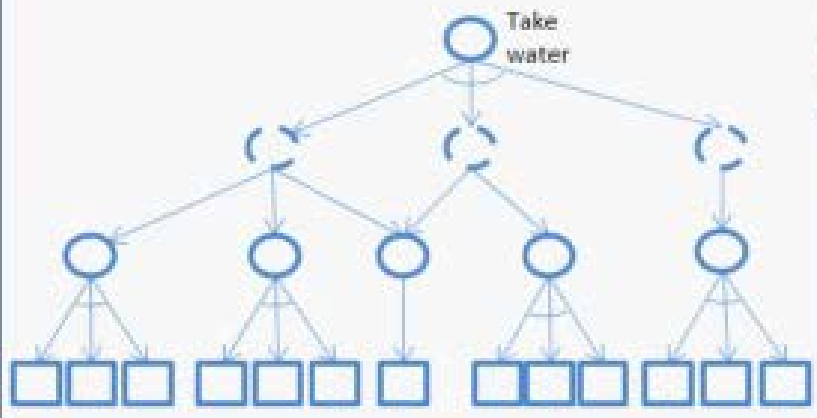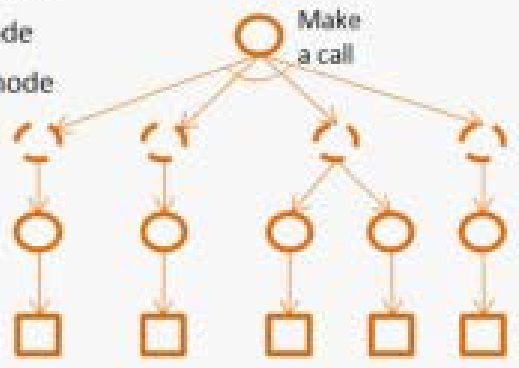# Handle event interruption



And-Node
Or-Node
Leaf-Node

First Partial parse tree of take water

Parse tree of take a phone

Second Partial parse tree of take water

Observed
Data

t

# Demo

And-node

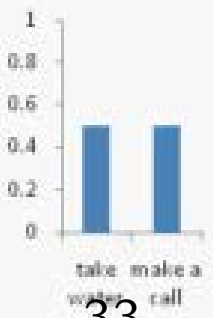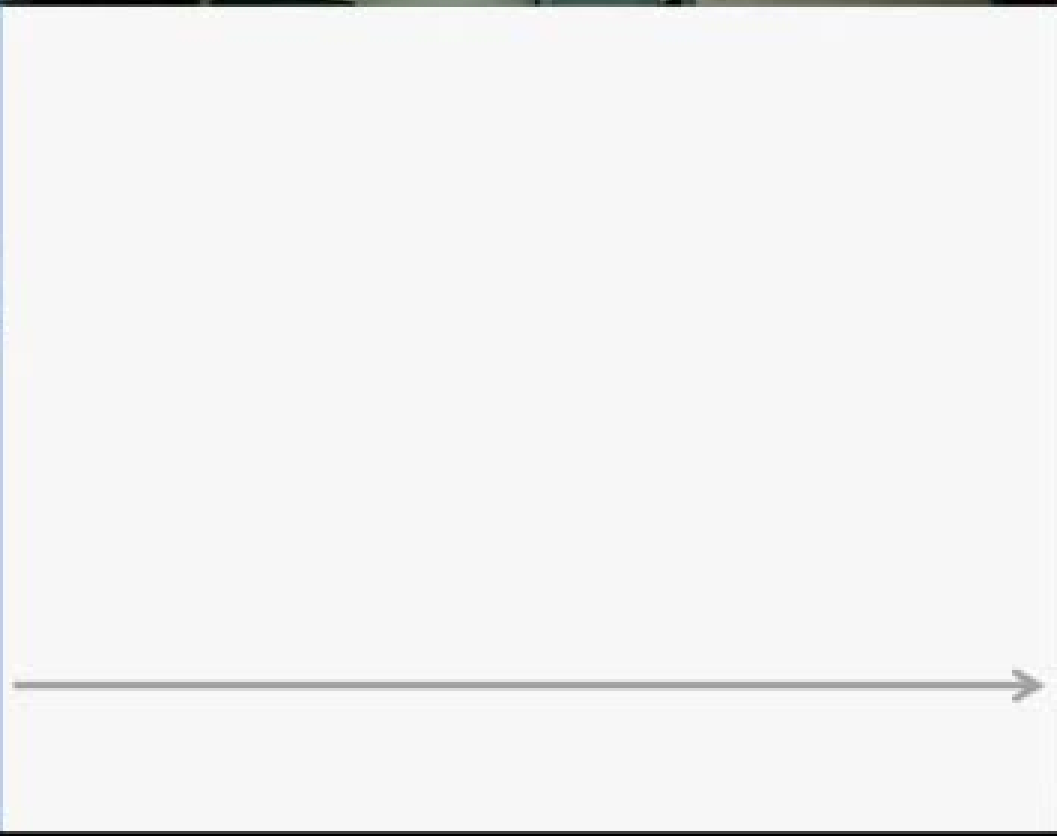Or-node

Leaf-node

Make a call

Take water

Intention

Prediction

take water    make a call

arrive trash can    arrive basin    arrive tea box    arrive wd    arrive phone
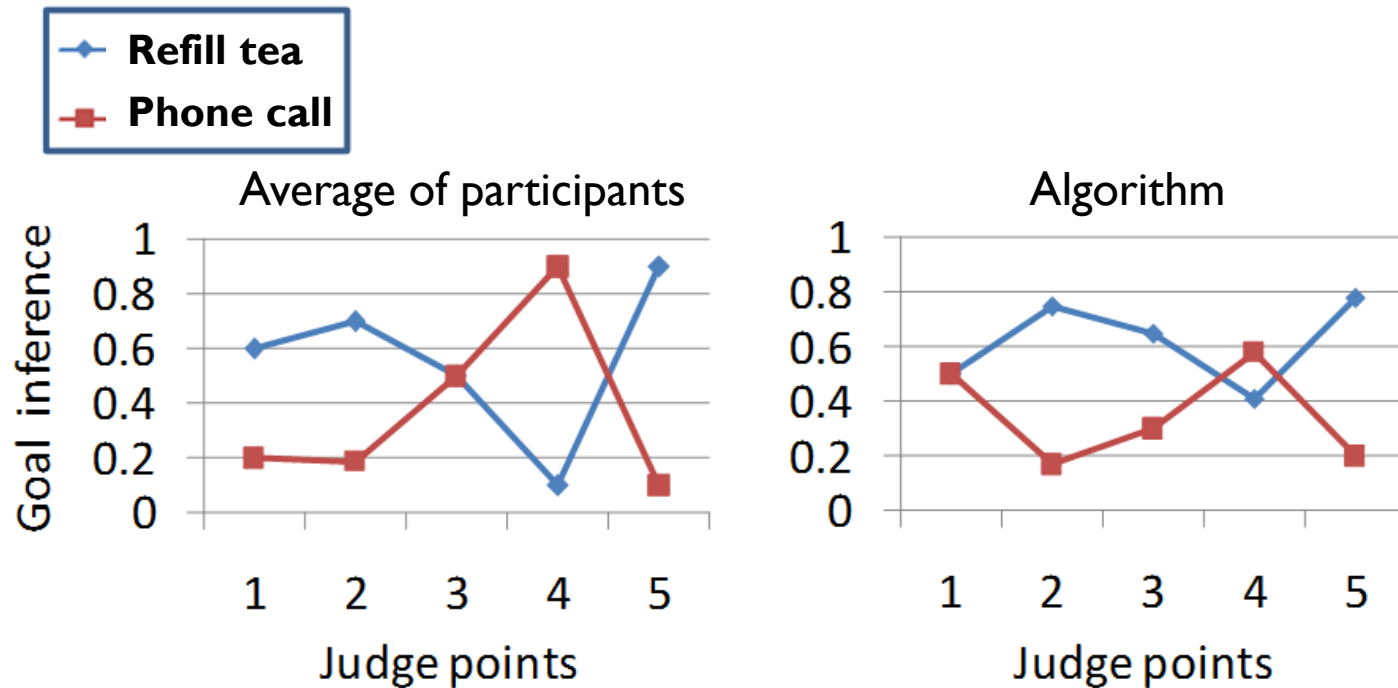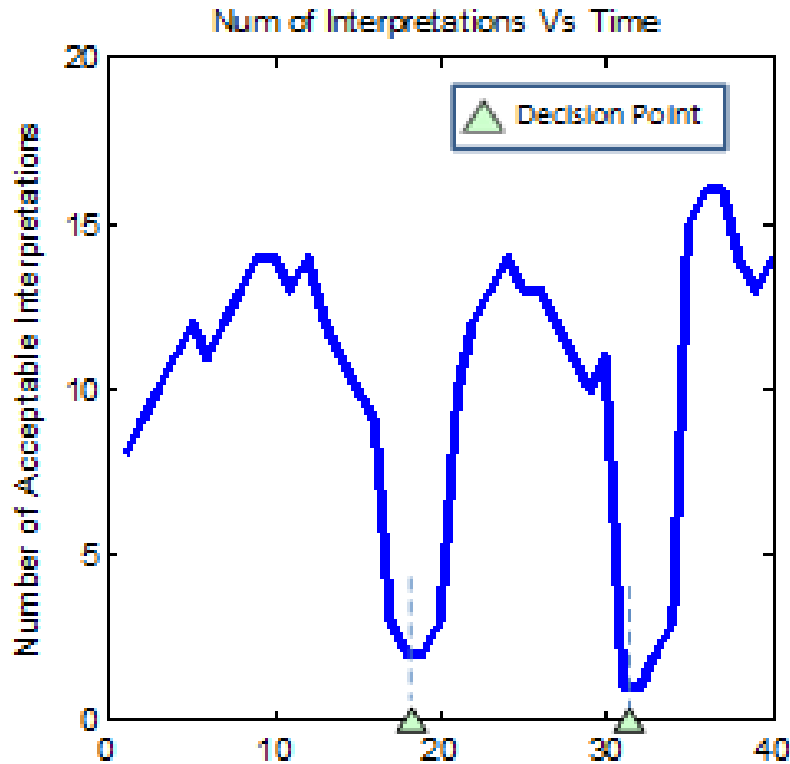
33

# Comparison with human prediction



M.T. Pei, Z.Z. Si, B. Yao, and S.C. Zhu, "Video Event Parsing and Learning with Goal and Intent Prediction," 2012

34

related work: [Baker, Saxe and Tenenbaum 2009]

# Computation complexity of parsing



Num of Interpretations Vs Time

Decision Point

- Initially the number of interpretations above a threshold grows rapidly over time.

- At certain decisive moments, i.e. when informative actions are observed, large number of unlikely interpretation drops below the threshold and hence is pruned.
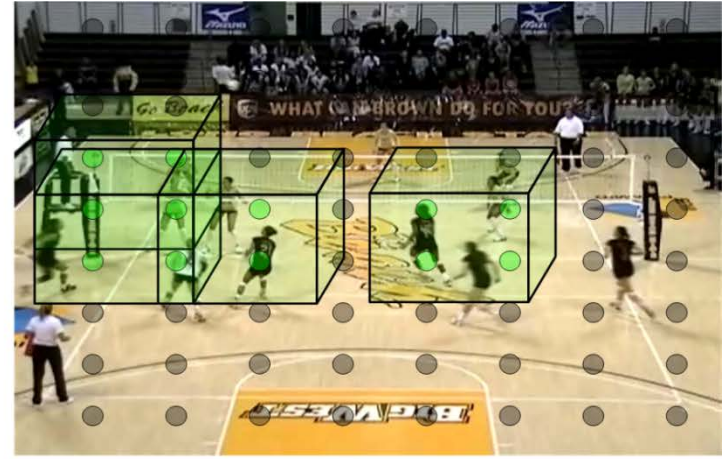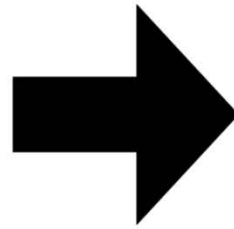
Pickup phone     Reach water boiler

# Weakly Supervised Learning of Temporal AND-OR Graph



Given Input Video  →  Classify & Localize

Stochastic activity has a random number of:

- actors,

- activity parts,

- spatiotemporal configurations

# Examples: Activities with Stochastic Structure

# Temporal AND-OR Graph

- AND nodes = Particular space-time configurations

- OR nodes = Alternative configurations

- Terminal nodes = BoWs

# Temporal AND-OR Graph

$$S(C) = 0.5(0.4x_1 P_1 + 0.2\overline{x_1}(1 - P_1) + \cdots$$
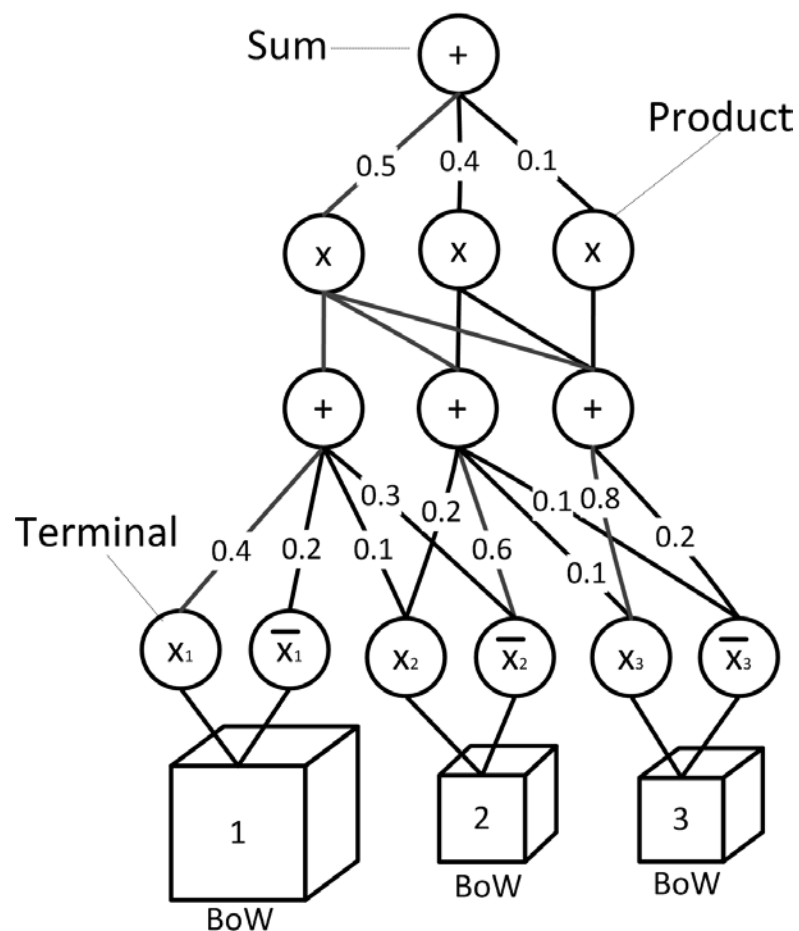
**Posterior:**

$$P(X|C) = S(C)/S_{X=1}$$

**OR nodes:**

$$S_i(C) = \sum_{j \in i^+} w_{ij} S_j(C)$$

**AND nodes:**

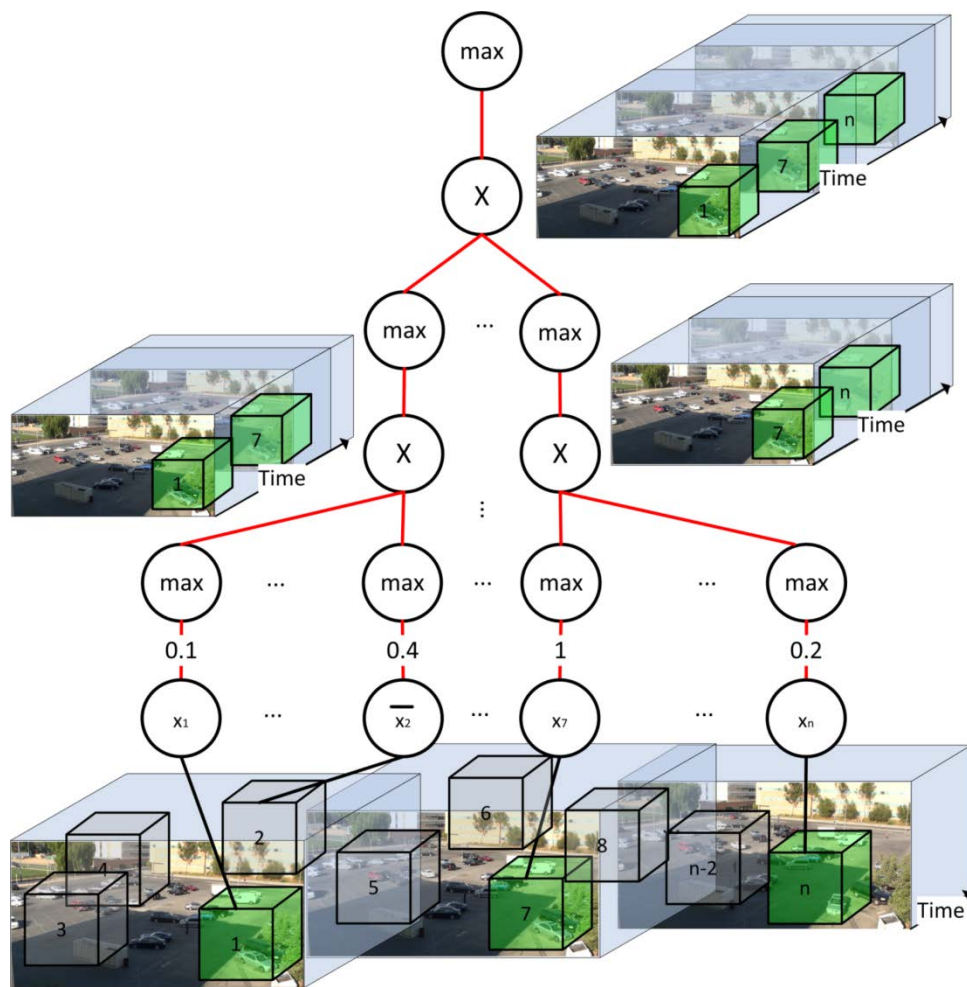$$S_k(C) = \prod_{l \in k^+} S_l(C)$$

# Learning – Variational EM

1. Learn AND-OR graph structure parameters – W
2. Learn Counting Grid parameters – $\pi$

$$V = \sum_t \Big[ \sum_b Q_b \log[(w_{ib1} x_b^t - w_{ib2} \overline{x_b^t})/Q_b]$$

$$+ \sum_b Q_b \sum_z (c_{bz}^t + \theta_z - 1) \log \Big[ \sum_{u \in H_b} \pi_{uz} \Big] \Big],$$

$$Q_b \propto \exp\Big[ \sum_{t,z} (w_{ib1} x_b^t - w_{ib2} \overline{x_b^t})(c_{bz}^t + \theta_z - 1) \log\big[ \sum_{u \in H_b} \pi_{uz} \big] \Big]$$

# Bottom up/Top Down Most Probable Explanation

$$\text{MPE: } \hat{a} = \text{argmax}_{a \in A} \hat{S}(C; a)$$

# Results – Volleyball Dataset

# Results –Volleyball Dataset