# Scheduling Top-Down and Bottom-Up Processes

Song-Chun Zhu, Sinisa Todorovic, and Ales Leonardis
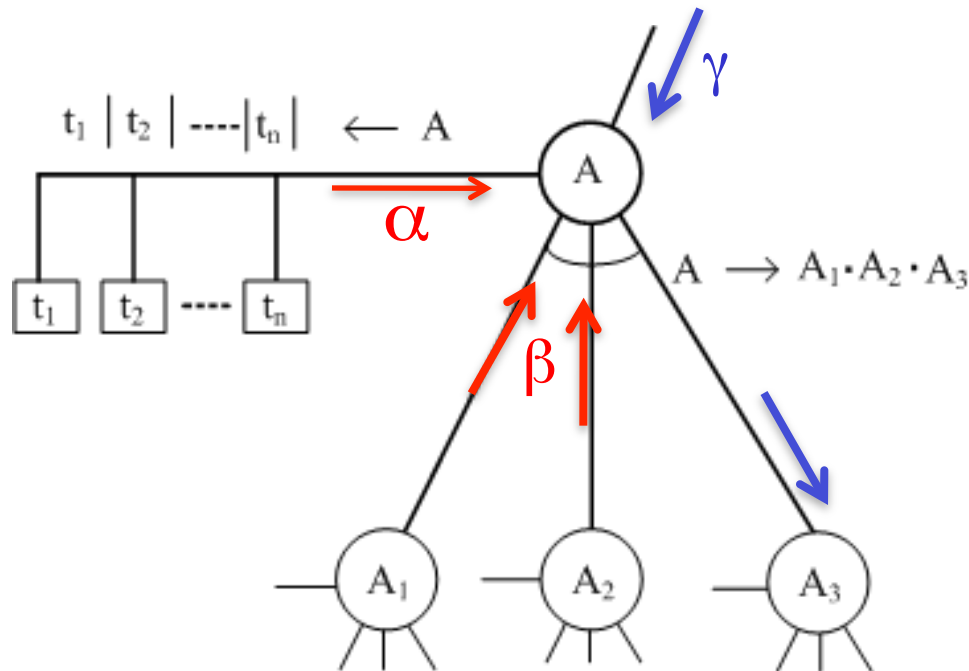
At CVPR, Providence, Rhode Island
June 16, 2012

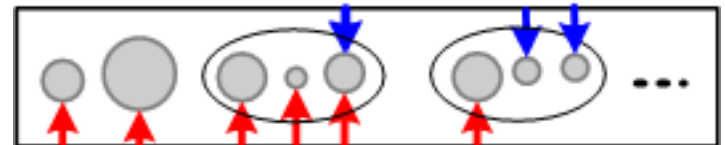# α, β and γ computing processes in AoG

The And-Or graph is a recursive structure.   So,  consider a node A.
  1.  Node A  at a coarse scale terminates to leaf nodes (ground)
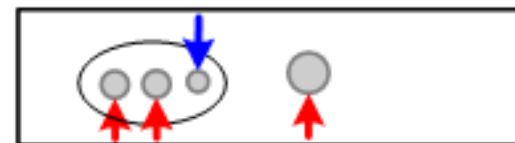  2.  Node A  is connected to the root.

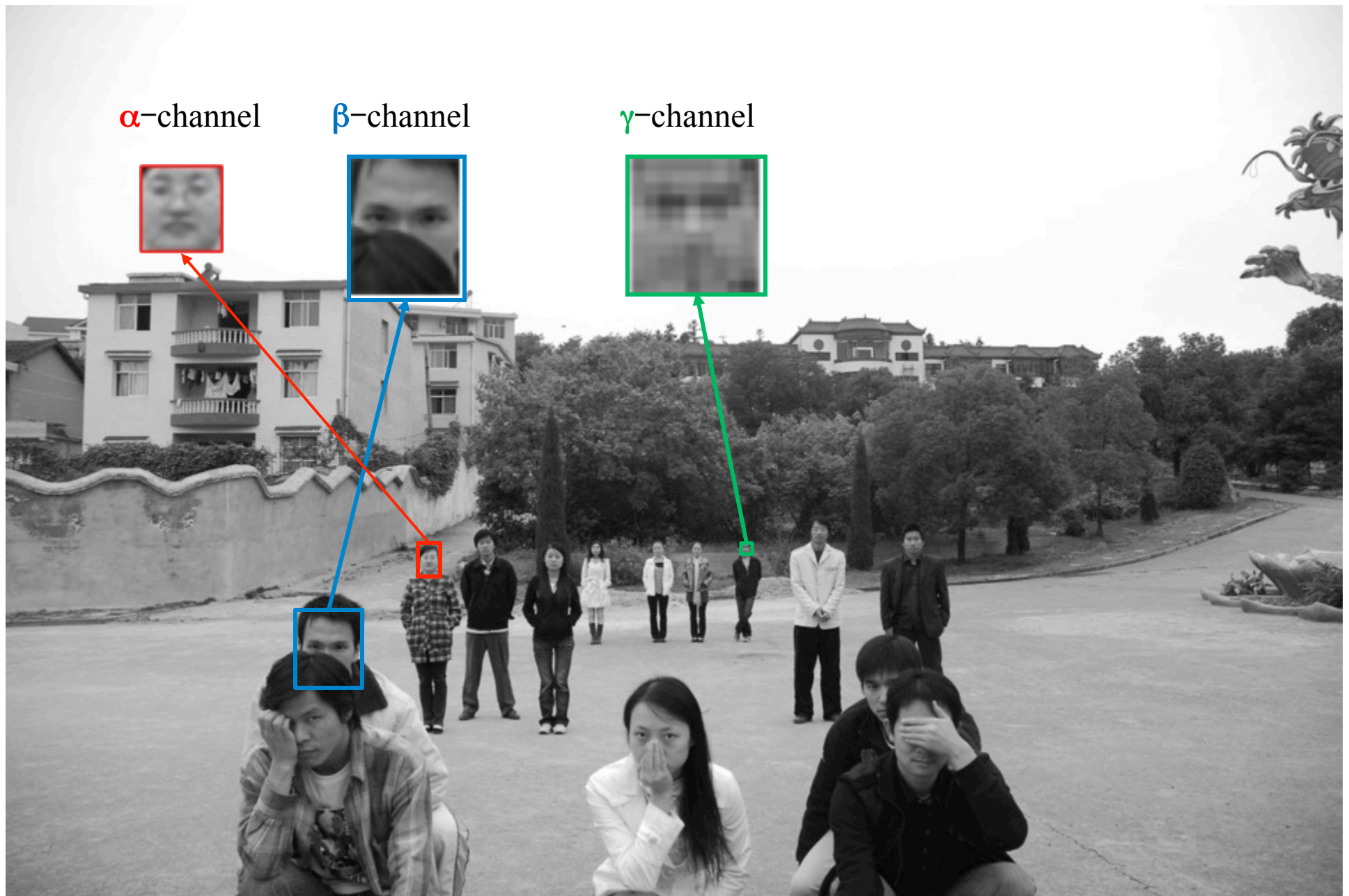Estimate the α/β/γ channels when they are applicable ---an optimal scheduling problem

# Example: Human Faces are Computed in 3 Channels

# Human faces in real scenarios

$\alpha$-channel

$\beta$-channel

$\gamma$-channel

# Hierarchical modeling and $\alpha$, $\beta$ and $\gamma$ computing

◯ And-node  ▢ terminate-node

Head-shoulder

$\gamma$: p(face | parents)

$\alpha$: p(face | compact image data)

Face

$\beta$: p(face | parts)

Left eye   Right eye   Nose   Mouse

1. Each node has its own $\alpha$, $\beta$ and $\gamma$ computing processes.

2. **How much does each channel contribute?**

Wu & Zhu IJCV11

# Hierarchical modeling and $\alpha$, $\beta$ and $\gamma$ computing



$$pg^{l*} = \arg\max_{pg} \left[ \log p(\wedge^l | \vee^l) + p(N) \underbrace{\sum_{i=1}^{N} \log p(X(\wedge_i^{l+}) | X(\wedge^l))} \right.$$

parse graph connectivity

$\gamma$: p(face | parents)

$\alpha$: p(face | image)

$\beta$: p(face | parts)

$$+ \log \underbrace{\frac{p(\Delta(t(\wedge^l)) | t(\wedge^l))}{q(\Delta(t(\wedge^l)))}}$$

$\alpha$ = detector

$$+ p(N) \Big[ \underbrace{\sum_{i=1}^{N} \log \frac{p(\Delta(t(\wedge_i^{l+})) | t(\wedge_i^{l+}))}{q(\Delta(t(\wedge_i^{l+})))} + \sum_{i \neq j} \log p(X(\wedge_i^{l+}), X(\wedge_j^{l+})) \Big]}$$

$\beta$ = zoom-in

$$+ \log \underbrace{\frac{p(\Delta(t(\wedge^{l-})) | t(\wedge^{l-}))}{q(\Delta(t(\wedge^{l-})))} + \log p(X(\wedge^l) | X(\wedge^{l-}))}$$

$\gamma$ = zoom-out

Wu & Zhu IJCV11

# $\alpha$ processes for the face node



$\alpha$-channels:  p(face | compact image data)

# β processes for the face node(when its α is off)



β-channels: p(face | parts), binding

α-channels of some parts are on
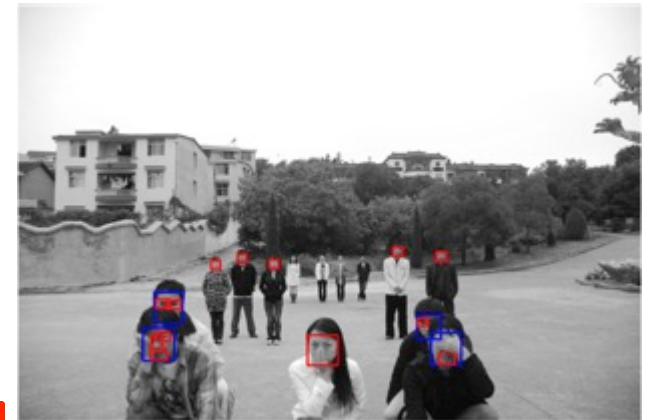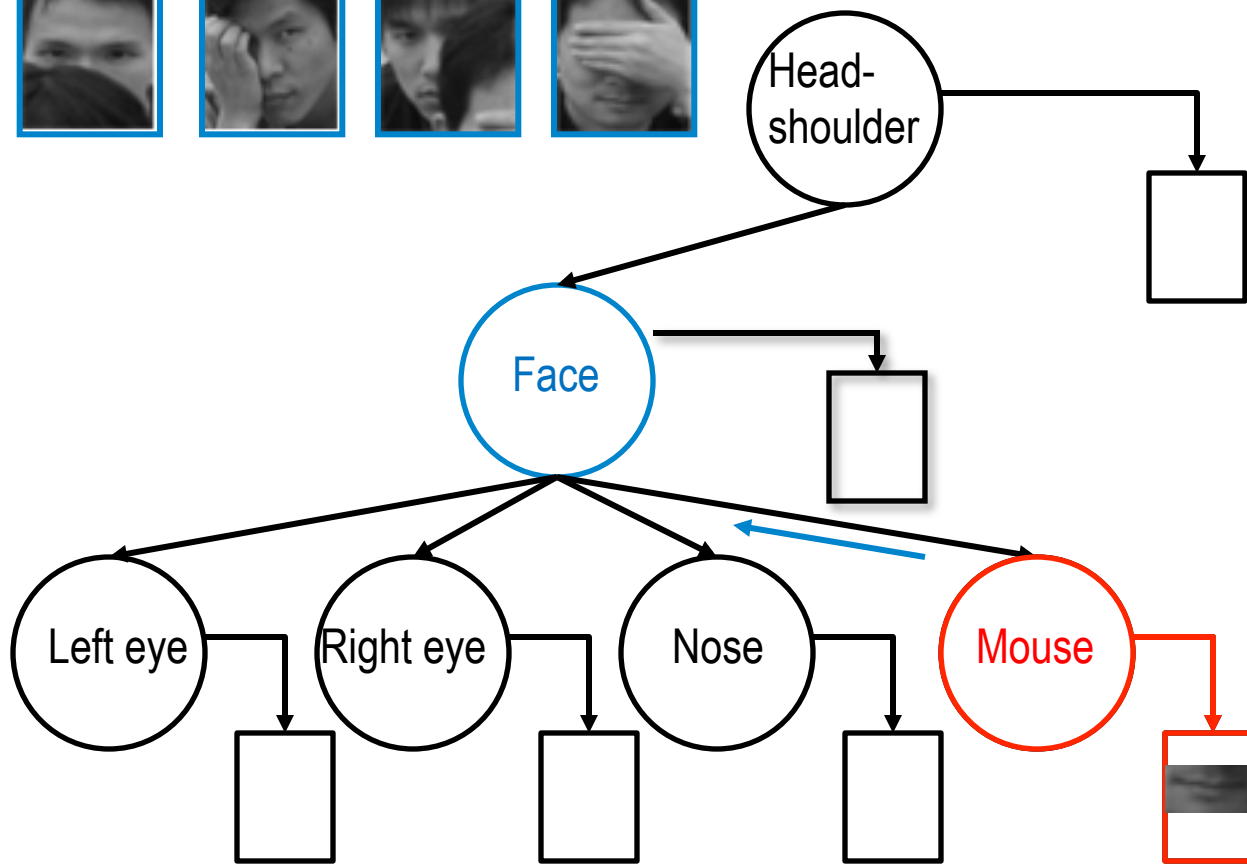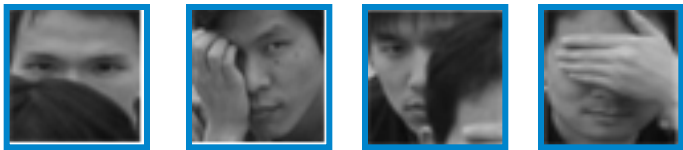
# β processes for the face node(when its α is off)

β-channels:  p(face | parts), binding



Head-shoulder
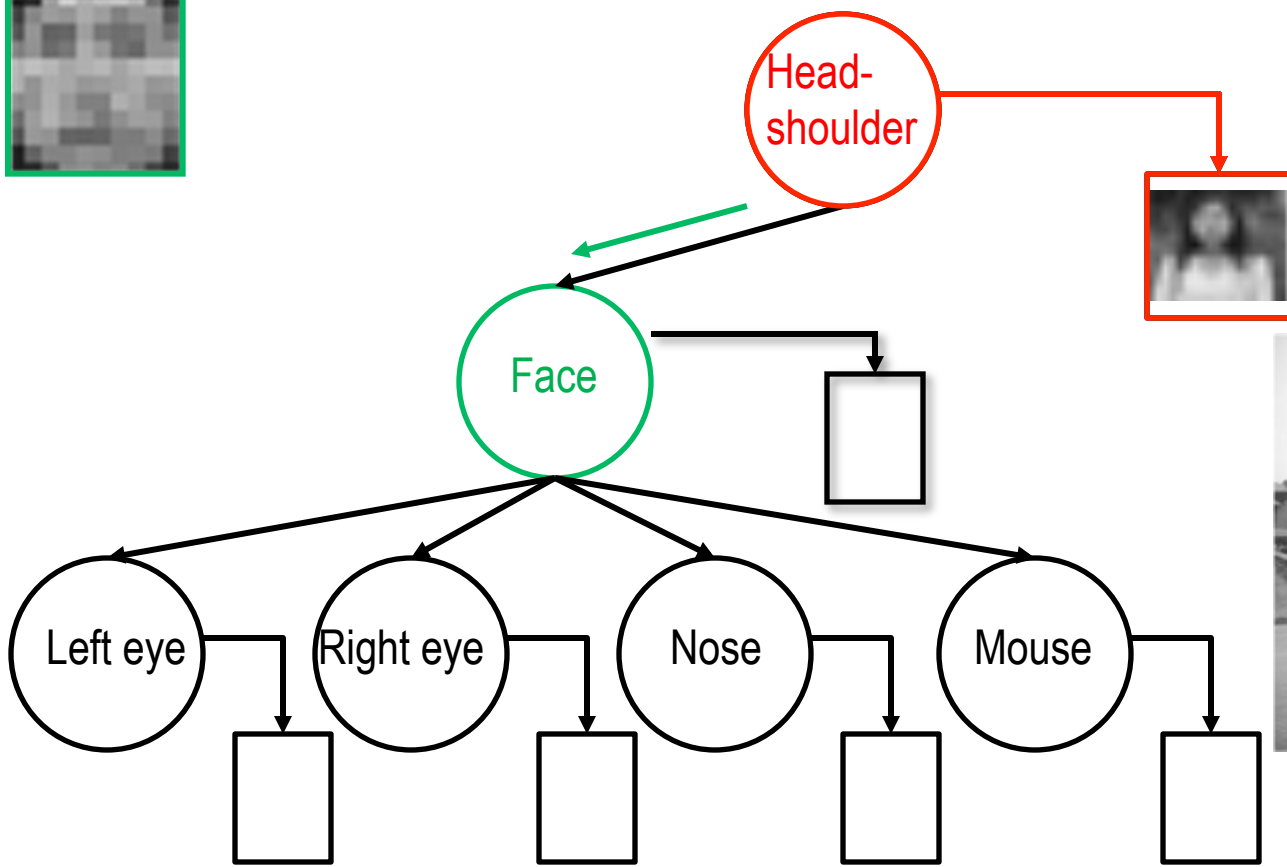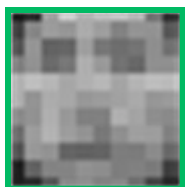
Face

Left eye    Right eye    Nose    Mouse

α-channels of some parts are on

# β processes for the face node(when its α is off)

β-channels:  p(face | parts), binding



α-channels of some parts are on

# β processes for the face node(when its α is off)

β-channels: p(face | parts), binding



α-channels of some parts are on

# γ processes for the face node(when it's α and β is off)

γ-channels:  p(face | parents), predicting

α-channels of some parents are on



Head-shoulder

Face

Left eye    Right eye    Nose    Mouse

# $\gamma$ processes for the face node(when it's $\alpha$ and $\beta$ is off)



$\gamma$-channels: p(face | parents), predicting

$\alpha$-channels of some parents are on

Head-shoulder

Face

Left eye

Right eye

Nose

Mouse

# γ processes for the face node(when it's α and β is off)

γ-channels: p(face | parents), predicting

α-channels of some parents are on



Head-shoulder

Face

Left eye    Right eye    Nose    Mouse

# γ processes for the face node(when it's α and β is off)

γ-channels:  p(face  | parents), predicting

α-channels of some parents are on



Head-shoulder

Face

Left eye

Right eye

Nose

Mouse

# In general: recursive $\alpha$, $\beta$ and $\gamma$ channels

# $\alpha$–channel: head-shoulder

# $\alpha$−channel: head-shoulder

# $\alpha$−channel: head-shoulder

# $\alpha$−channel: face

# $\alpha$–channel: face

# $\alpha$–channel: face

# $\alpha$-channel: eye

# $\alpha$−channel: eye

# $\alpha$-channel: eye

# α−channel: nose

# $\alpha$–channel: nose

# $\alpha$−channel: nose

# $\alpha$–channel: mouth

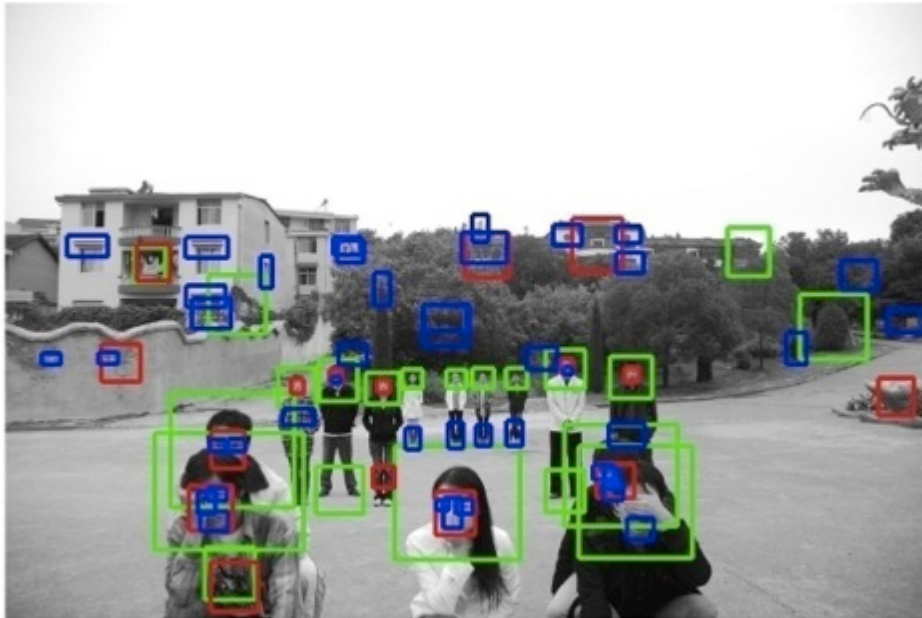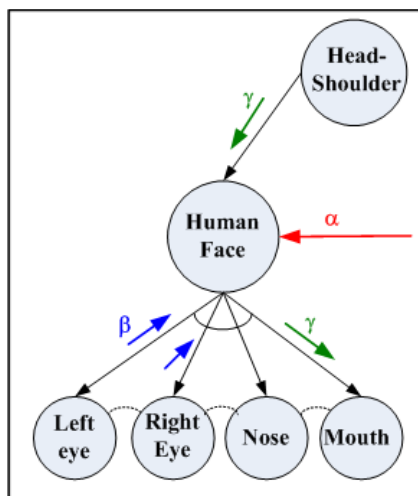# $\alpha$–channel: mouth

# $\alpha$−channel: mouth

# All α channels
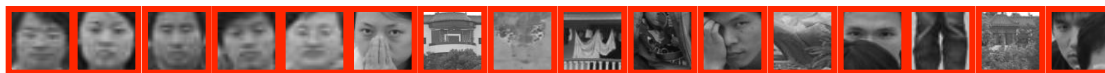


keep things and throw away stuffs by

integrating α, β and γ channels

# Integrating α, β and γ channels

# Integrating $\alpha$, $\beta$ and $\gamma$ channels

# Integrating α, β and γ channels

# Integrating α, β and γ channels

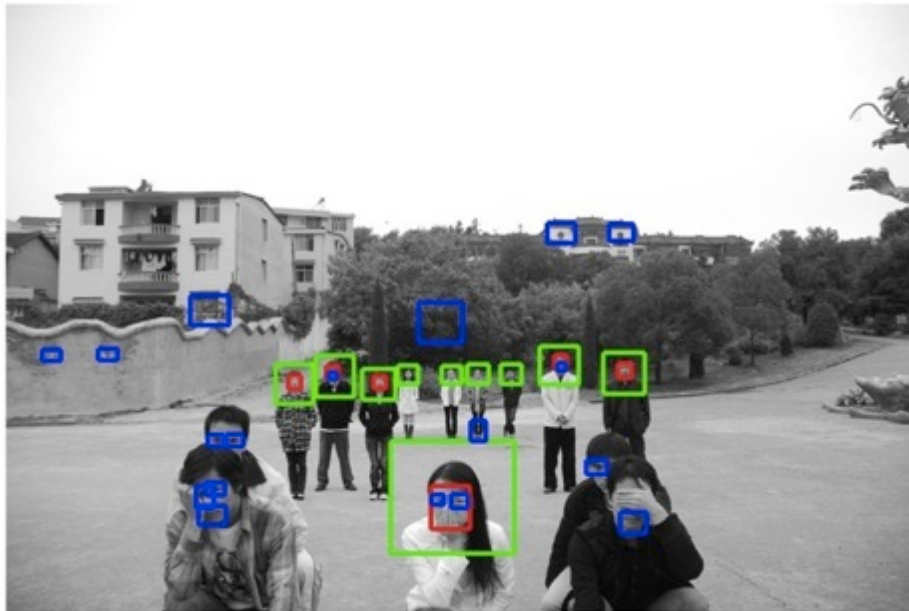# Integrating α, β and γ channels

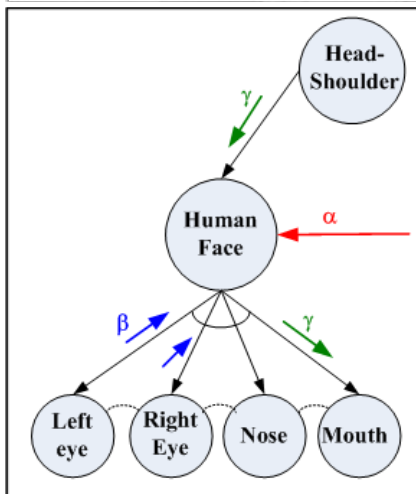# Integrating α, β and γ channels

# Integrating α, β and γ channels

# Integrating α, β and γ channels

# Integrating α, β and γ channels

# Integrating α, β and γ channels

# Integrating $\alpha$, $\beta$ and $\gamma$ channels

# Integrating $\alpha$, $\beta$ and $\gamma$ channels

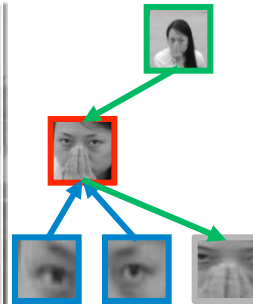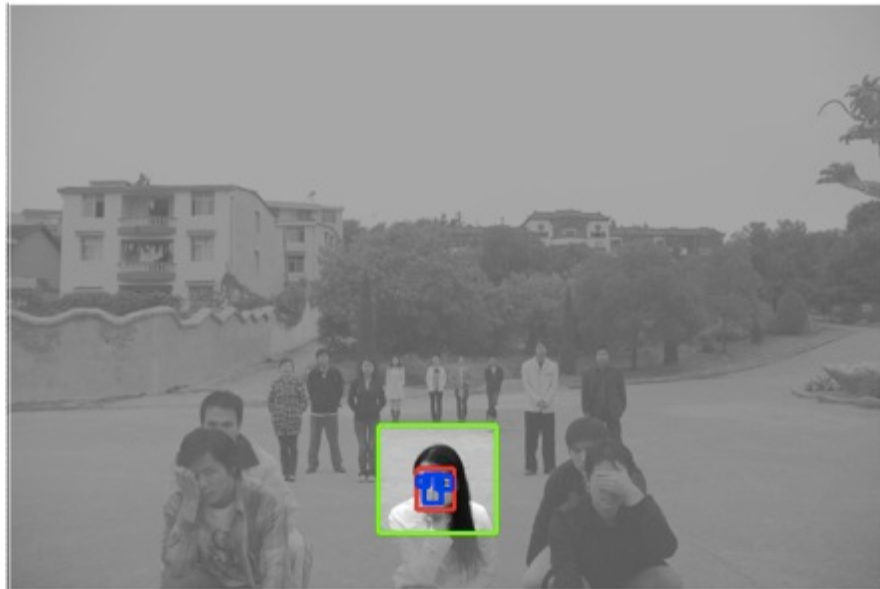# Integrating $\alpha$, $\beta$ and $\gamma$ channels

# Integrating $\alpha$, $\beta$ and $\gamma$ channels

# Integrating α, β and γ channels
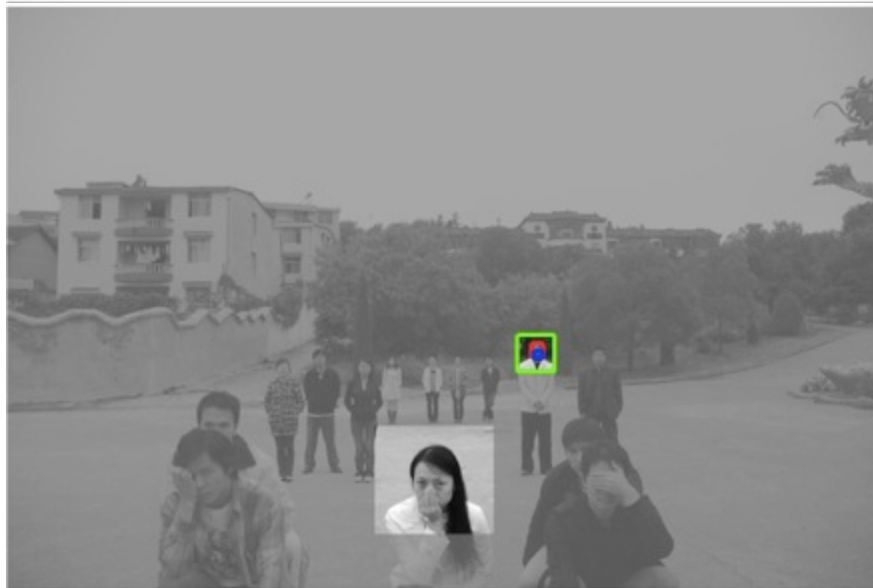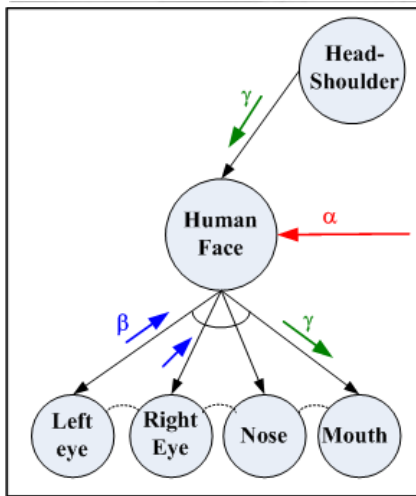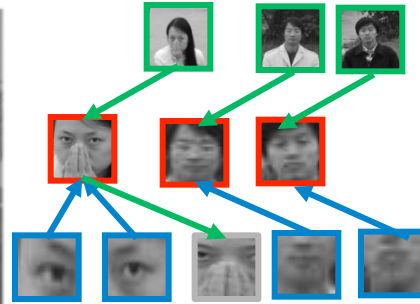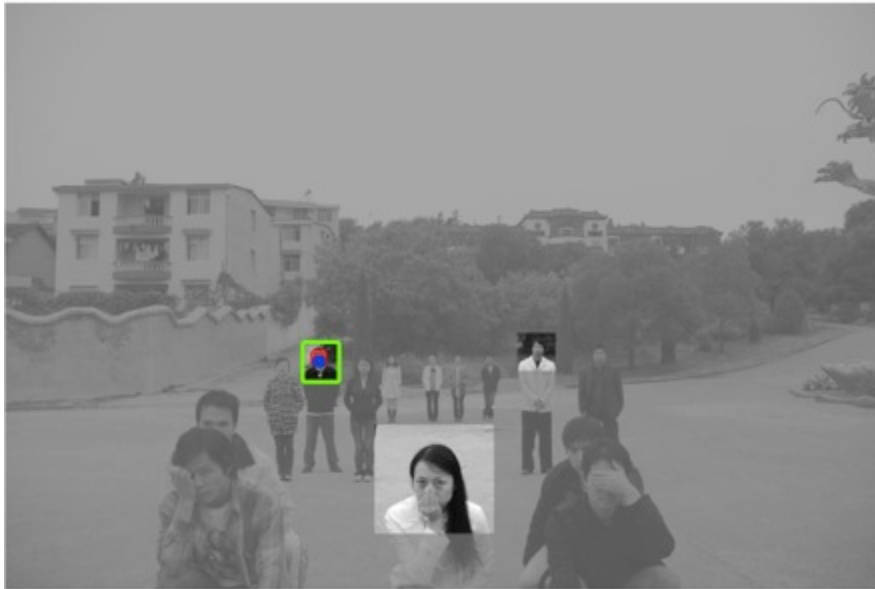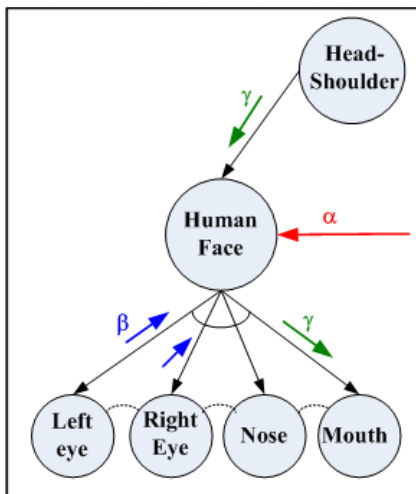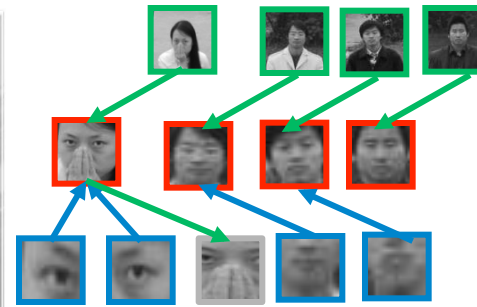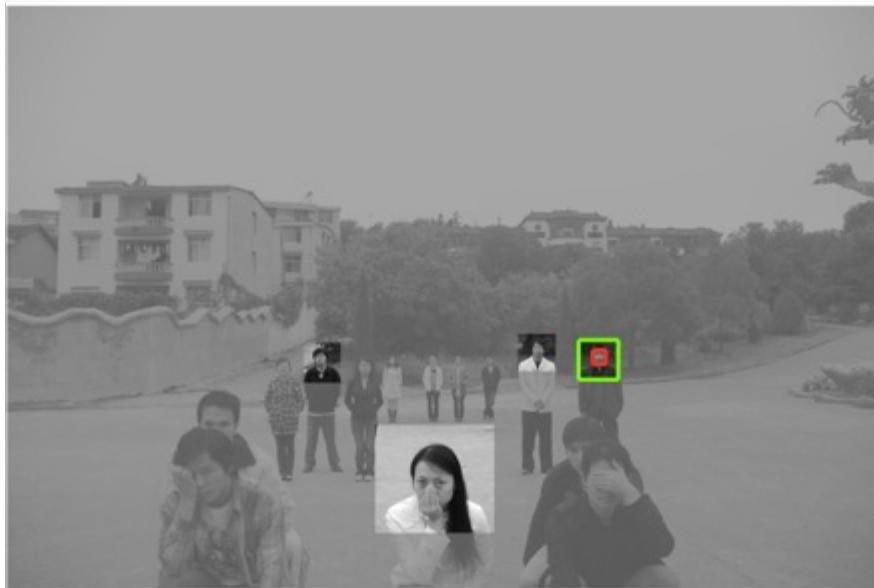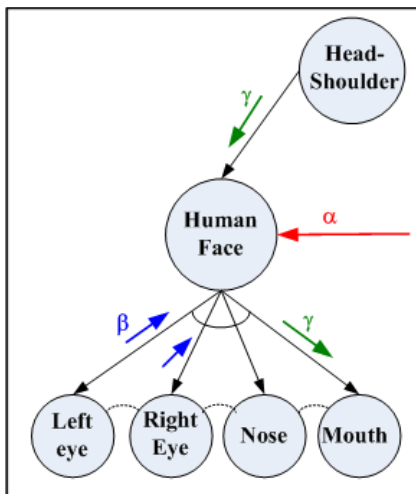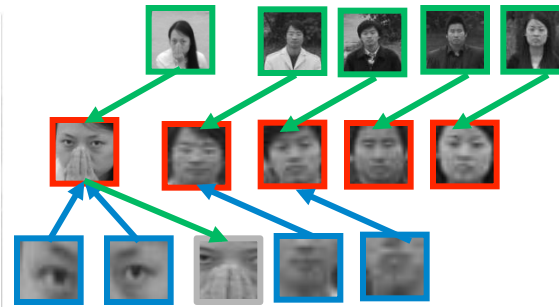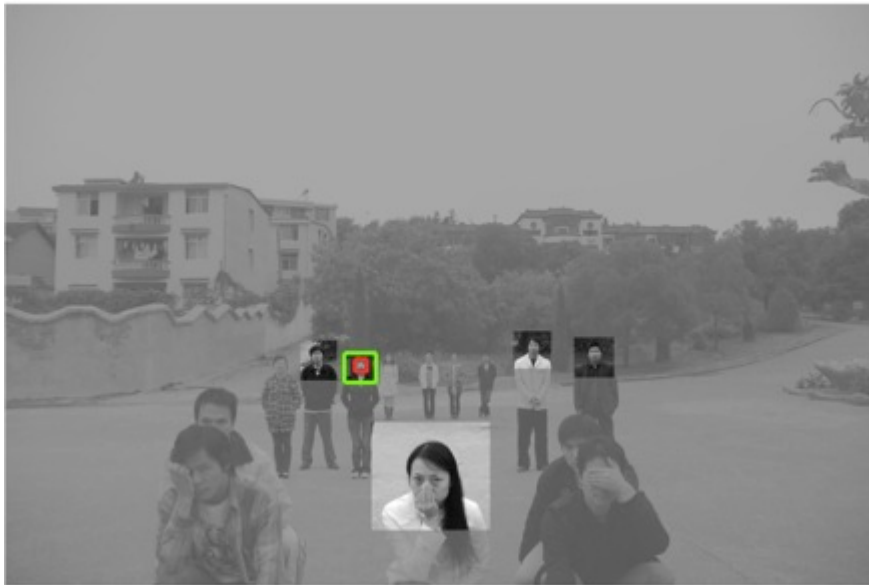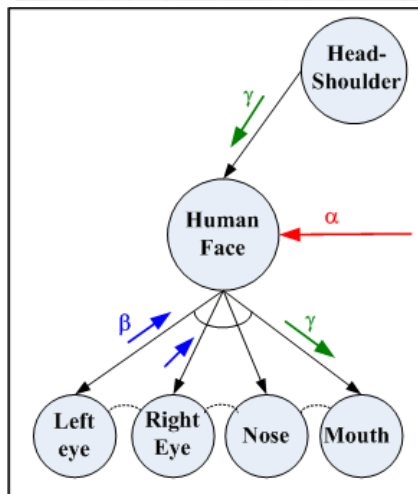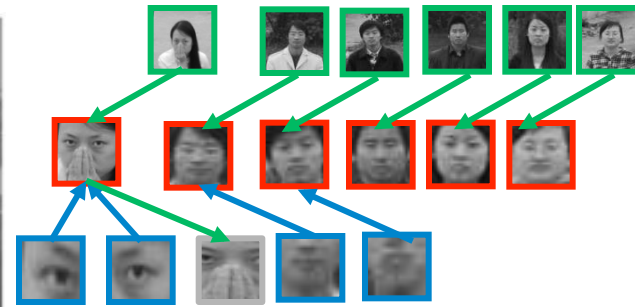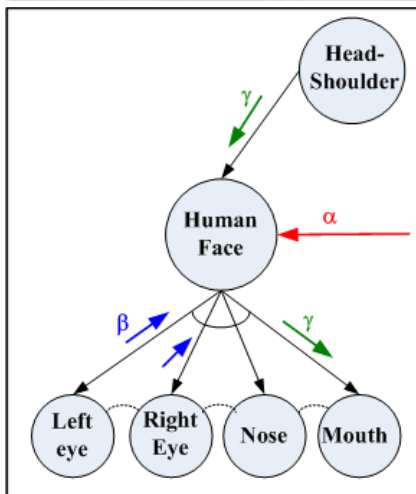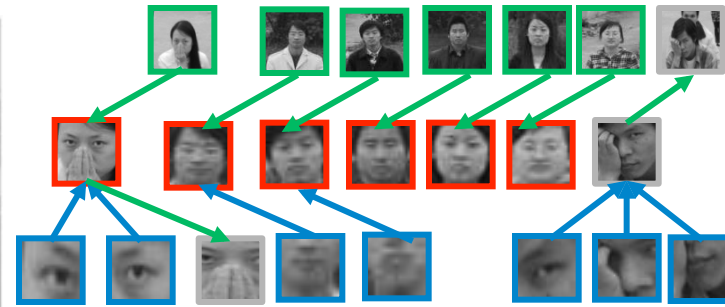
# Integrating α, β and γ channels

# Information contribution

$\alpha$ channels

# Performance improvement

red for α, blue for α+β, green for α+γ, cyan for α+β+γ channels

# Cost-Sensitive and Goal-Guided Inference in AoG

Song-Chun Zhu, Sinisa Todorovic, and Ales Leonardis

At CVPR, Providence, Rhode Island
June 16, 2012

# Problem



Given a high-resolution and long video, showing

- Spatially large scenes,

- Many people engaged in a wide range of

  - individual actions

  - group activities

# Problem



Goal:

Answer WHAT, WHERE, and WHEN queries about

- individual actions
- group activities

# Multiscale Event Understanding

# A Principled Framework Needed for

- **Goal-driven** zooming-in or zooming-out,

- **Cost-sensitive** inference via

  – Scheduling of the bottom-up and top-down inference

  – Any-time inference -- adapt to a given time budget

# And-Or Graph and Parse Graph



$$pg^{l*} = \arg\max_{pg}\left[\log p(\wedge^l|\vee^l) + p(N)\sum_{i=1}^{N}\log p(X(\wedge_i^{l+})|X(\wedge^l))\right.$$

parse graph connectivity

$$+ \log\frac{p(\Delta(t(\wedge^l))|t(\wedge^l))}{q(\Delta(t(\wedge^l)))}$$

$\alpha$ = detector

$$+ p(N)\left[\sum_{i=1}^{N}\log\frac{p(\Delta(t(\wedge_i^{l+}))|t(\wedge_i^{l+}))}{q(\Delta(t(\wedge_i^{l+})))} + \sum_{i\neq j}\log p(X(\wedge_i^{l+}), X(\wedge_j^{l+}))\right]$$

$\beta$ = zoom-in

$$+ \log\frac{p(\Delta(t(\wedge^{l-}))|t(\wedge^{l-}))}{q(\Delta(t(\wedge^{l-}))} + \log p(X(\wedge^l)|X(\wedge^{l-}))$$

$\gamma$ = zoom-out

# Inference -- Alpha of Individual Action



Waiting

# Inference -- Alpha of Individual Action



Waiting    Walking

# Inference -- Alpha of Individual Action



Waiting  Walking

# Inference -- Bottom Up (Beta)



Group Walking   Group Discuss

# Inference – Alpha of Group Activities



Group Walking  Group Discuss

# Inference -- Top Down (Gamma)



Group Walking    Group Discuss

# Inference -- Top Down (Gamma)

# Inference -- Alpha of Objects



Stairs  Chairs  Food bus

# Inference -- Bottom Up (Beta)



Waiting  Walking

# Inference



Waiting  Walking

# Queries about Individual Actions



Only Alpha

Zoom-out

Zoom-in and Zoom-out

# Scheduling

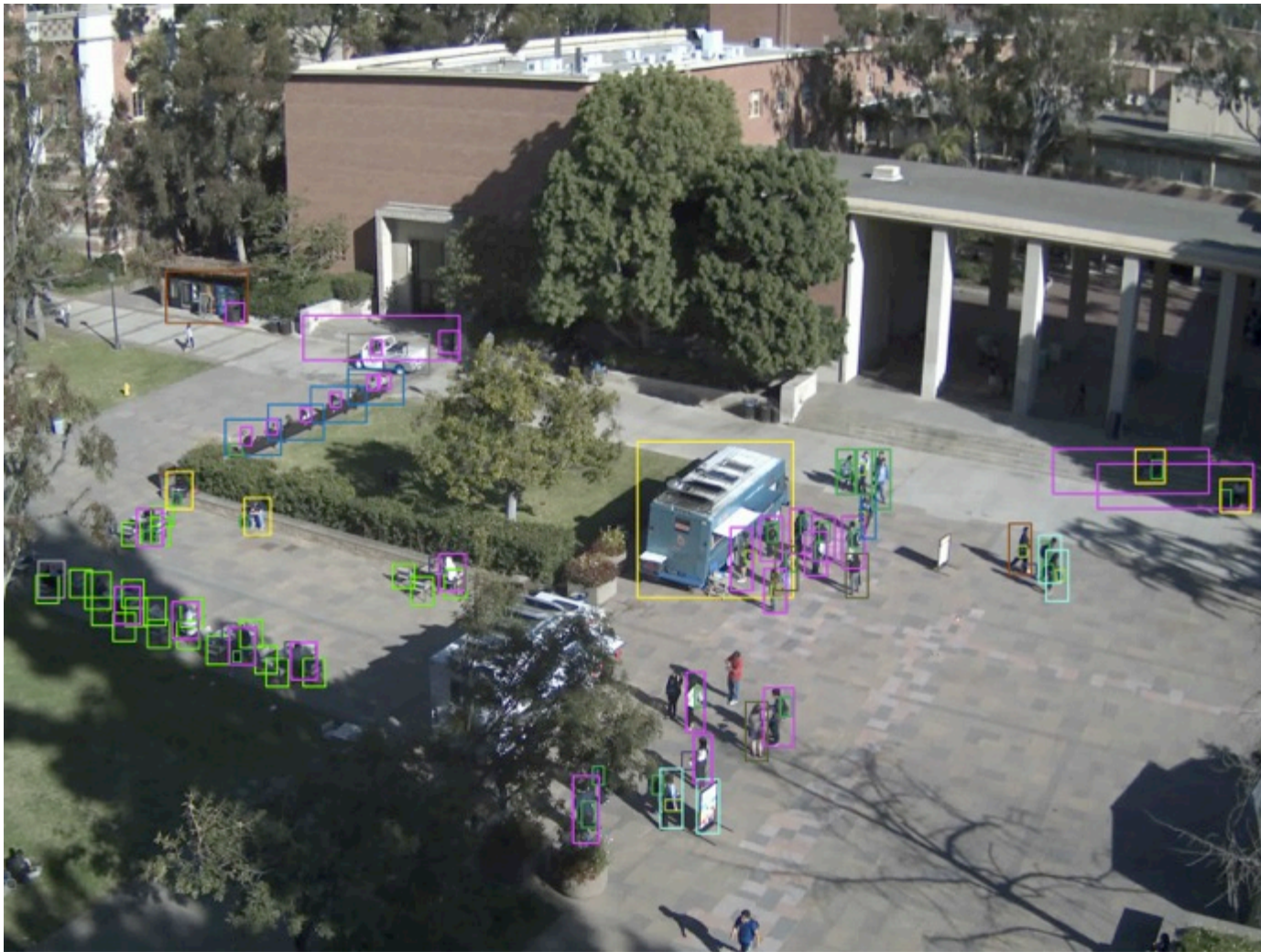$$\mathbb{V}_{\pi_{\mathcal{Q}}}^{(\tau)}(s) = \mathbb{E}[G(s,a)] + \sum_{s'} \mathbb{T}(s, \pi_{\mathcal{Q}}(s, \tau), s') \cdot \mathbb{V}_{\pi_{\mathcal{Q}}}^{(\tau-1)}(s'), \quad \tau = 1, ..., \mathbb{B}$$

- States (s) =
  - Query,
  - All previous actions,
  - Their rewards until this point
- Actions (a) =
  - Run detector and terminate (Alpha)
  - Bottom-up (Beta)
  - Top-down (Gamma)

# Policies, Given a Query (Q)

1. Alpha(Q)

2. Alpha(Q), Bottom-up, Alpha(Context), Top-down

3. Alpha(Q), Top-down, Alpha(Details), Bottom-up

4. Alpha(Q), Bottom-up, Alpha(Context), Top-down,
   Alpha(Details), Bottom-up

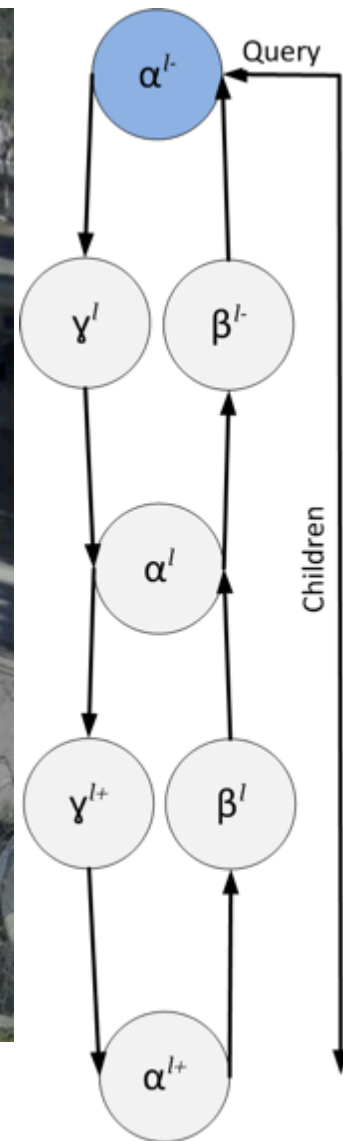5. Alpha(Q), Top-down, Alpha(Details), Bottom-up
   Alpha(Context), Top-down
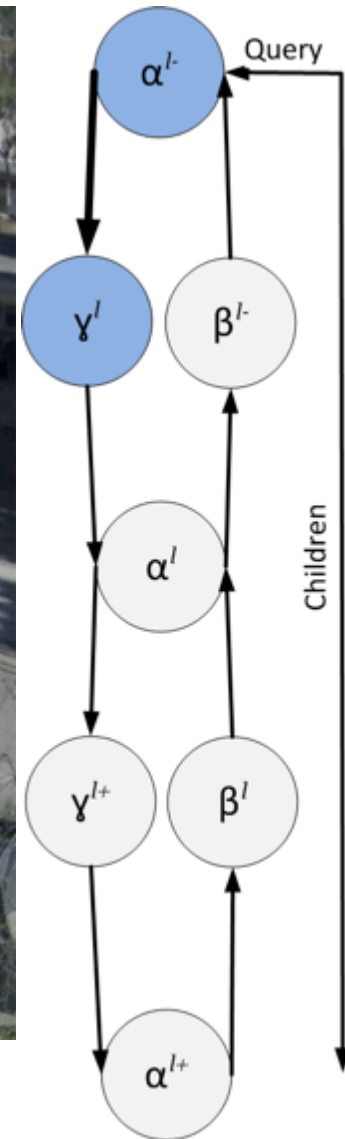
# Group Query

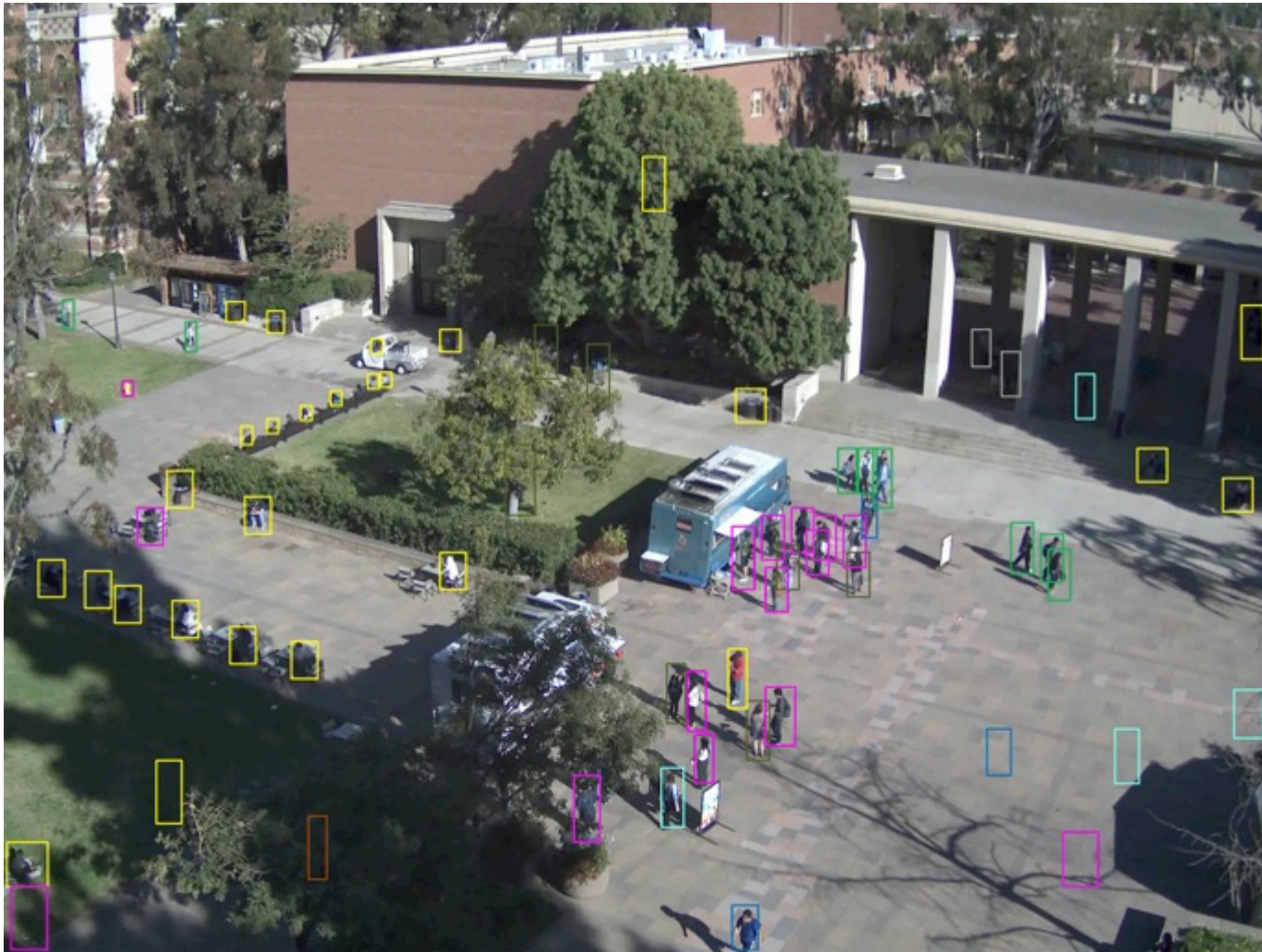# Group Query Policy #1
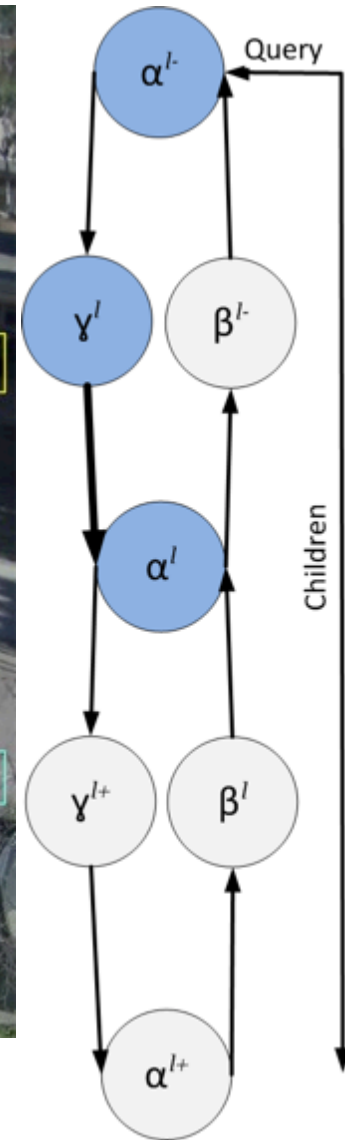


Group Walking    Group Discuss

# Group Query Policy #1



Group Walking ■ Group Discuss ■ ■ ■

# Group Query Policy #1



Waiting    Walking

# Group Query Policy #1



Group Walking     Group Discuss

# Group Query Policy #2



Group Walking    Group Discuss

# Group Query Policy #2



Group Walking   Group Discuss

# Group Query Policy #2

# Group Query Policy #2
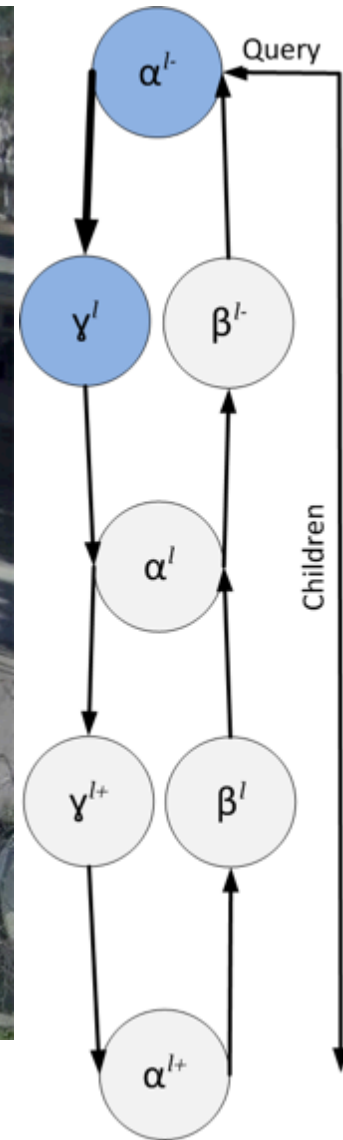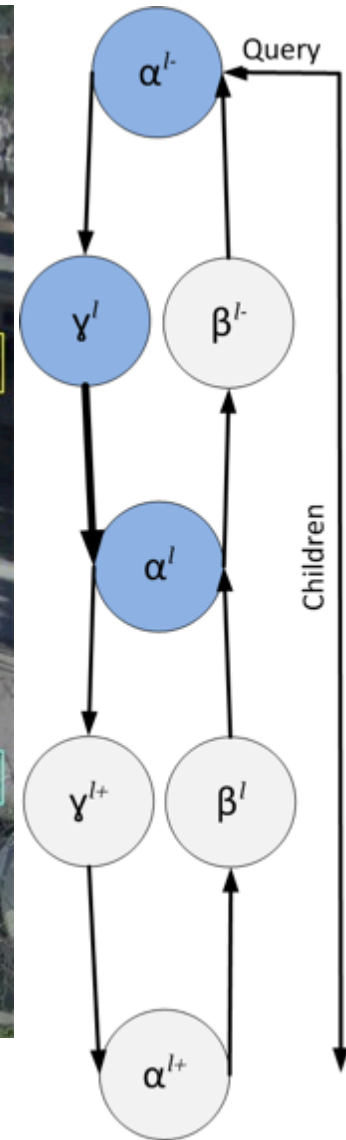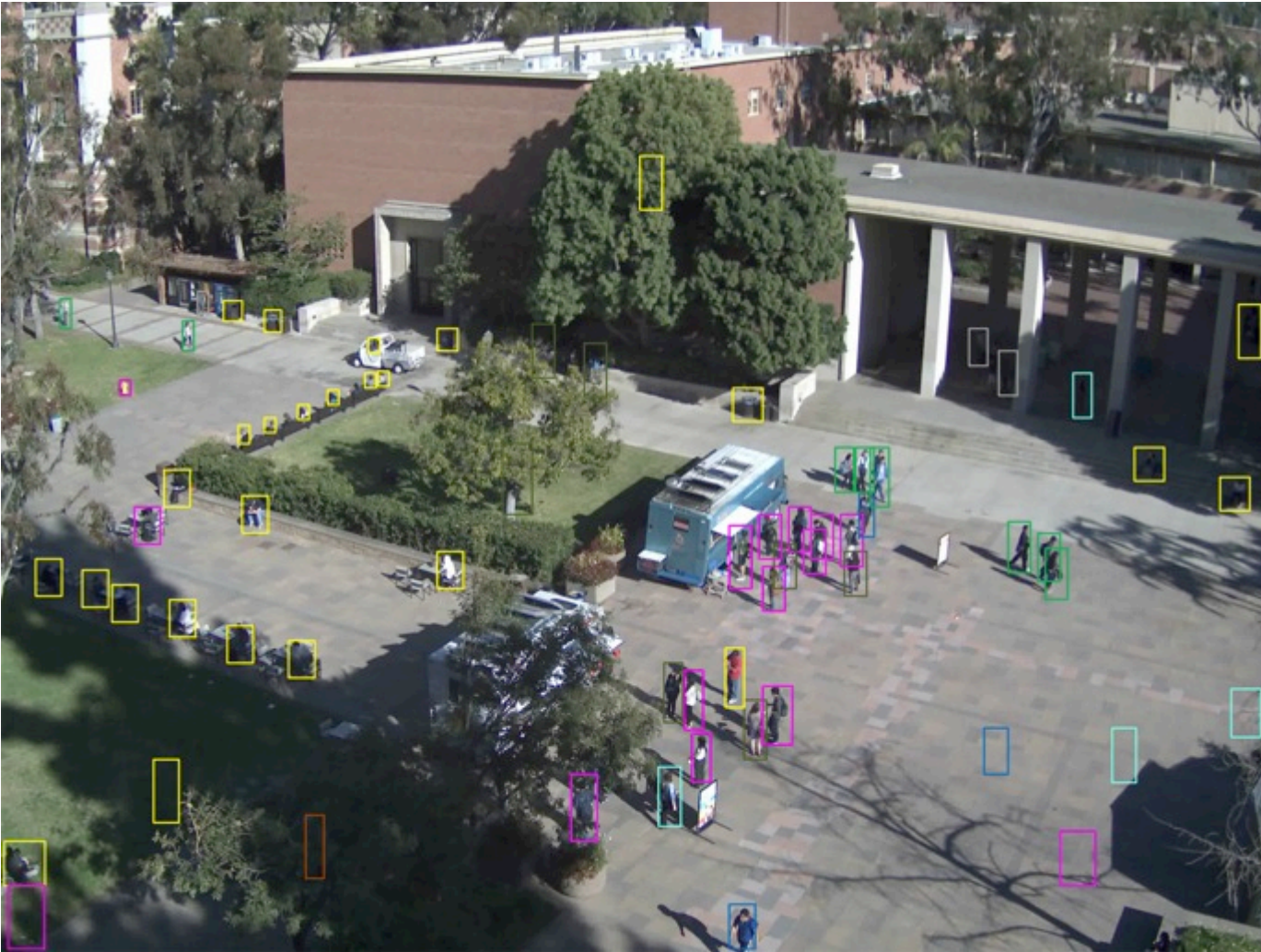
# Group Query Policy #2

# Group Query Policy #2

# Group Query Policy #2



Group Walking    Group Discuss

# Group Query Policy #2
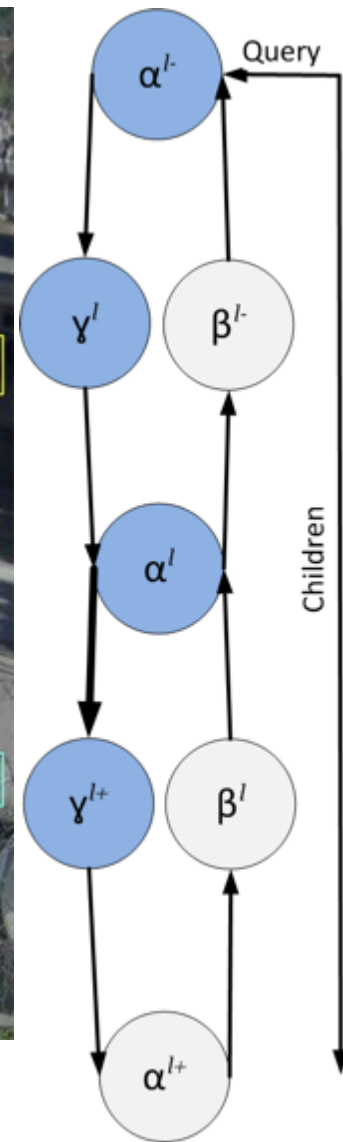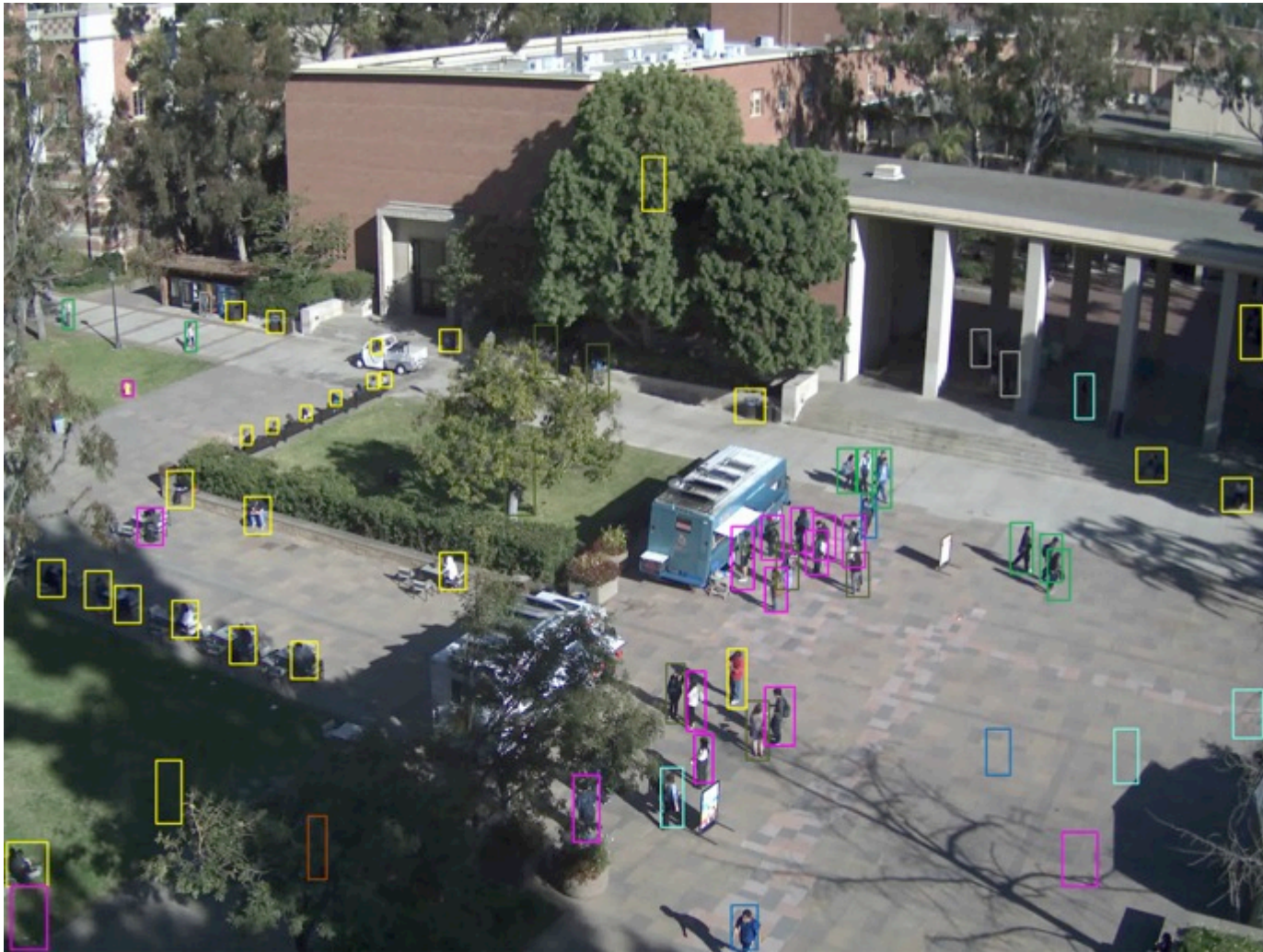


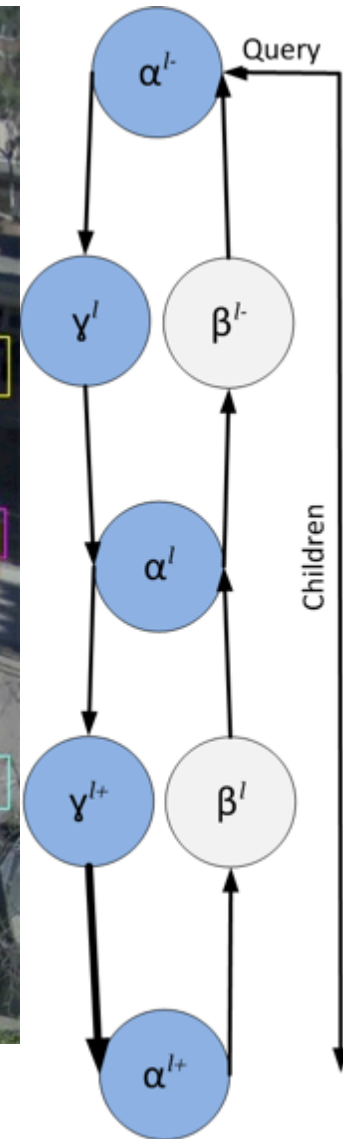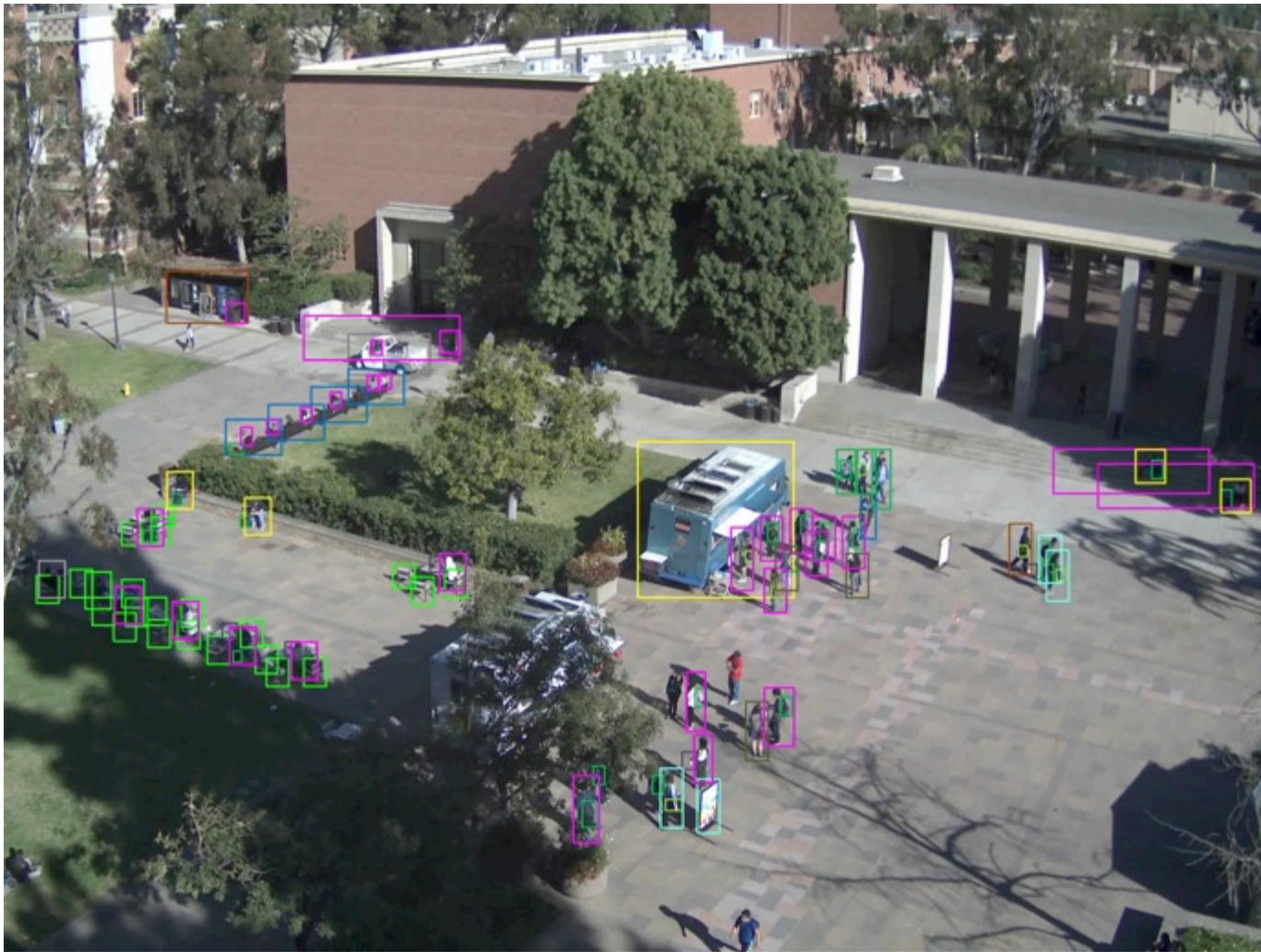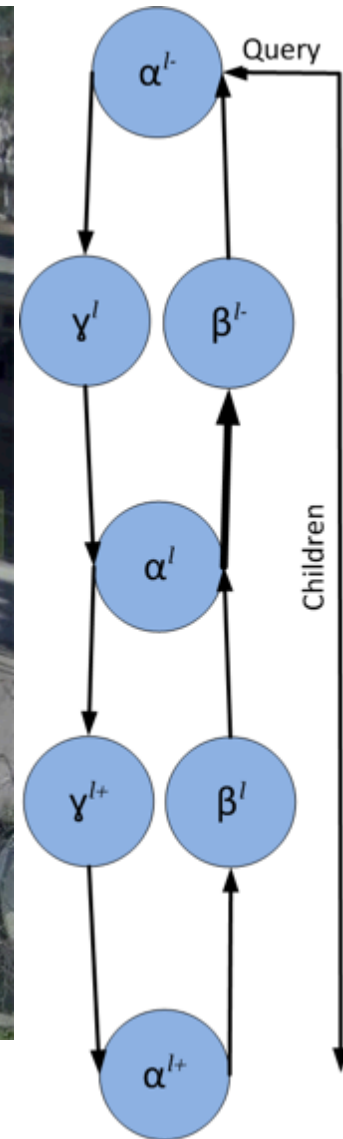Group Walking    Group Discuss
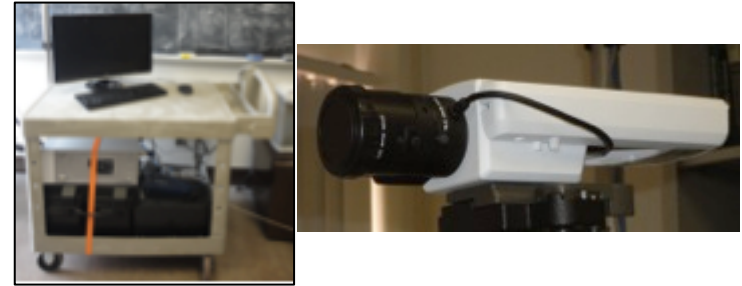
# Group Query



Alpha

Policy #1

Policy #2

# New Dataset



- Footage: 106min
- Frames: 12 fps
- Res.: 2560x1920 pixels
- Portable acquisition system



**VIRAT**

Benchmark datasets:
- single actor
- single group
- single action

# Domain Knowledge

- 5 Group Activities:
  - Walking together, Queuing, Campus tour, ...

- 10 Individual Actions:
  - Walking, Sitting, Riding a bike, ...

- 16 Objects:
  - Food truck, Vending machine, Bike, Backpack, ...