

---

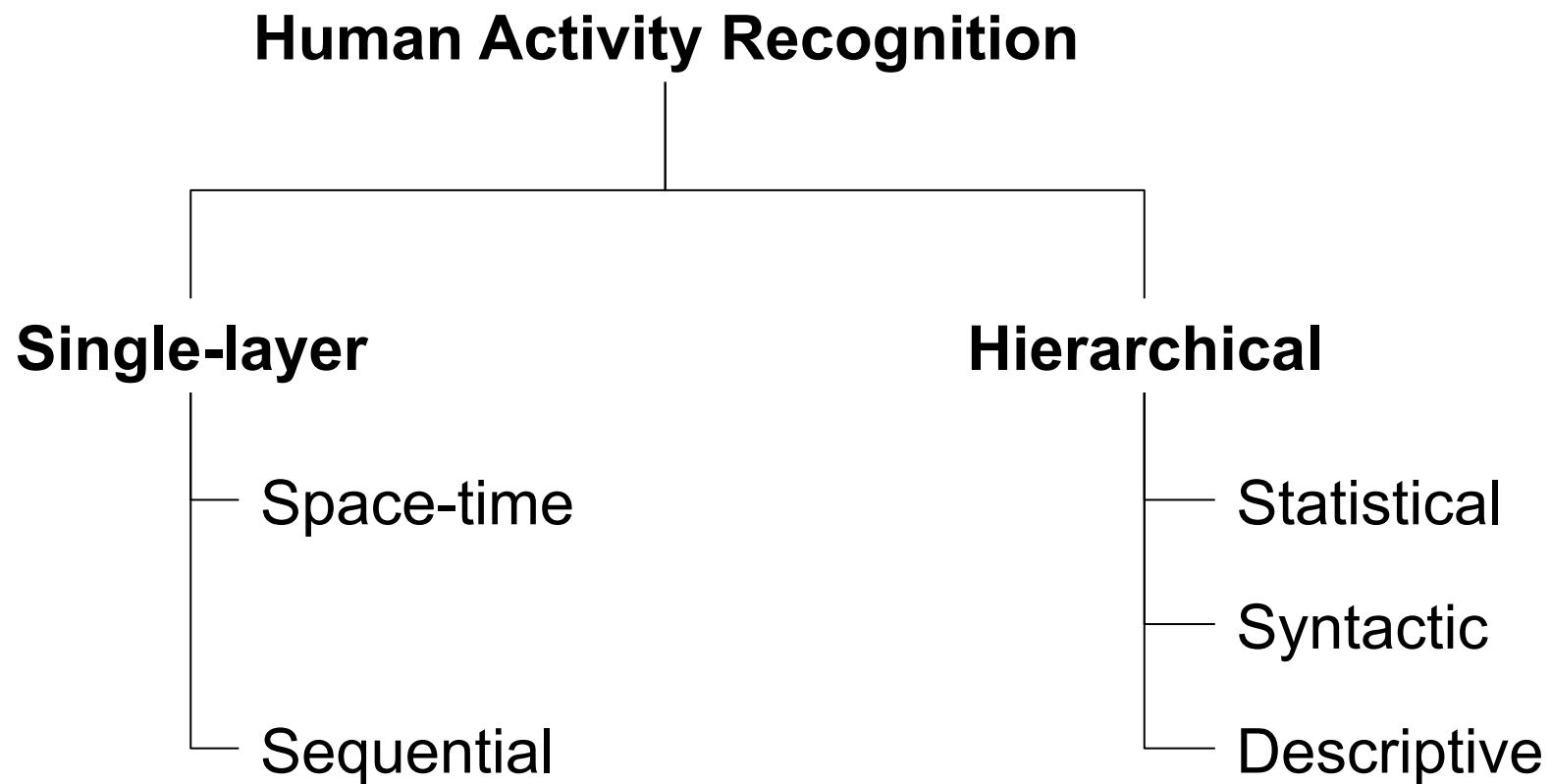
# Frontiers of *Human Activity Analysis*

J. K. Aggarwal  
Michael S. Ryoo  
Kris M. Kitani



# Overview

---



# Motivation

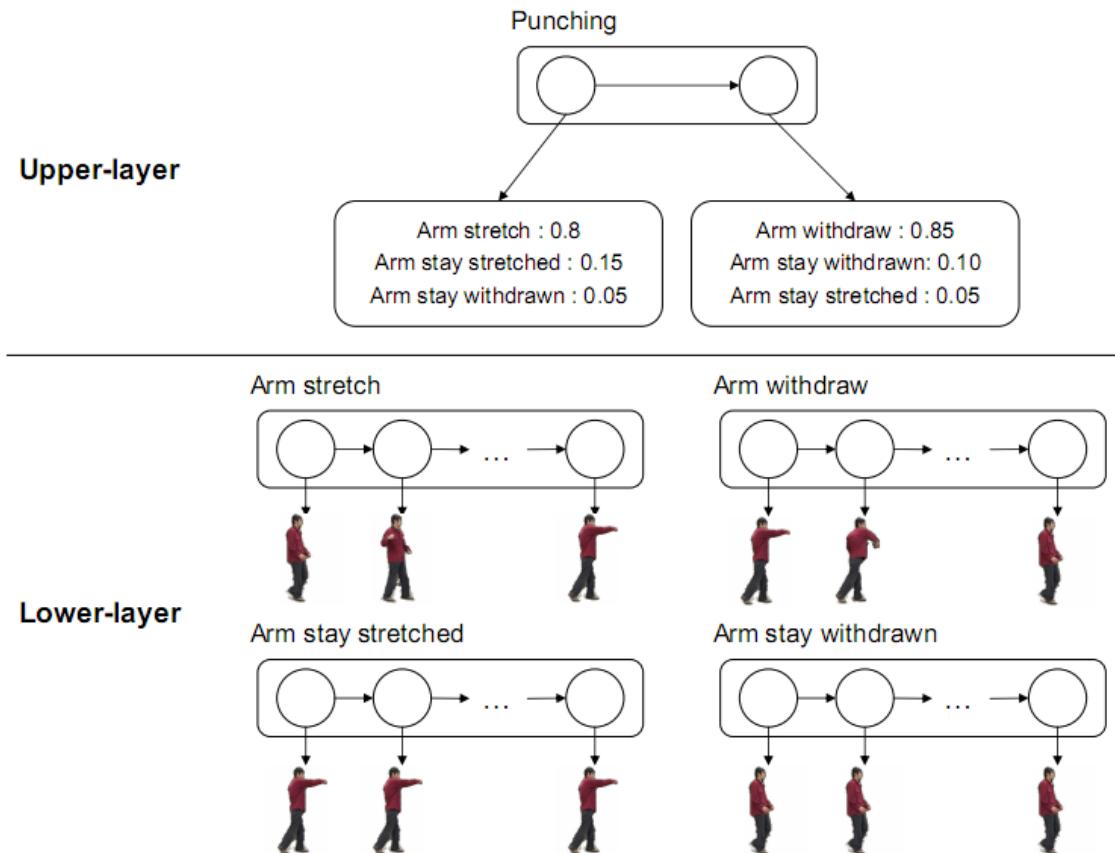
---

How do we interpret a sequence of actions?



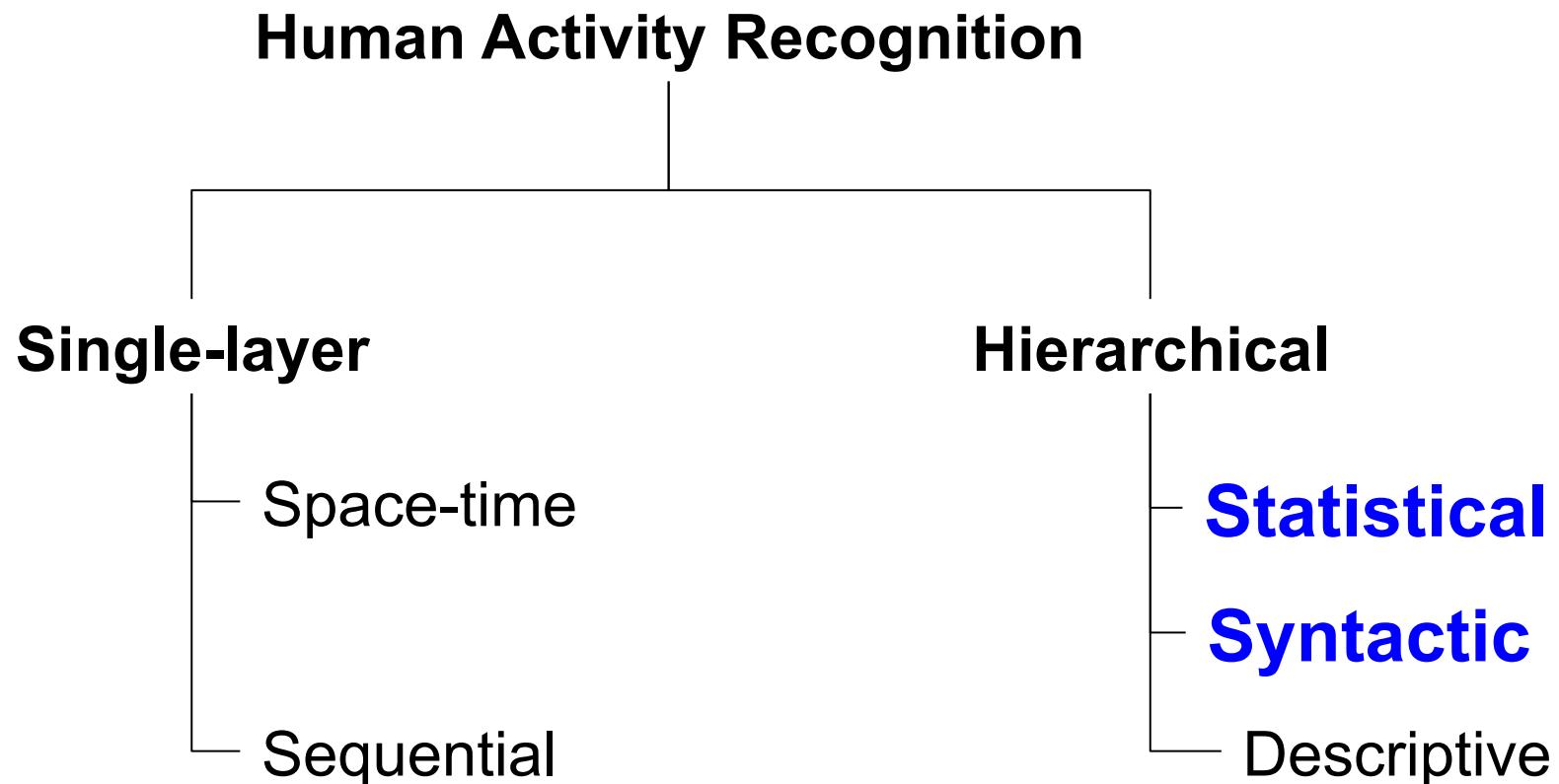
# Hierarchy

Hierarchy implies decomposition into sub-parts



# Now we'll cover...

---



---

# Syntactic Approaches

# Syntactic Models

---

Activities as strings of symbols.

s t r i n g s o f s y m b o l s

What is the underlying structure?

# Early applications to Vision

Tsai and Fu 1980.

Attributed Grammar-A Tool for Combining Syntactic and Statistical Approaches to Pattern Recognition.

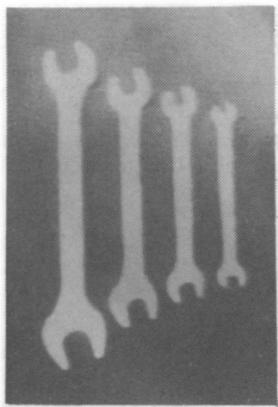
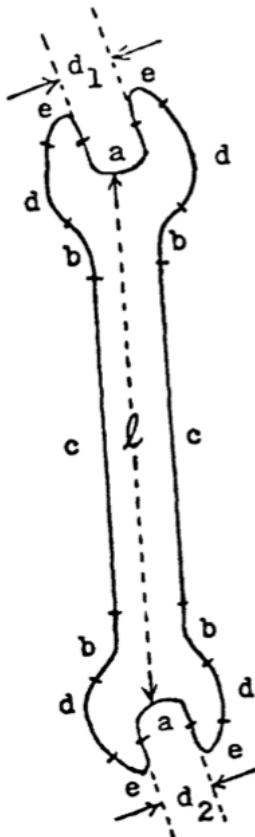


Fig. 6. Set of I wrenches.



$$z_I = aedbcbdeaedbcbde$$

Fig. 8. I wrench and its boundary primitives.

# Hierarchical syntactic approach

---

- Useful for activities with:
  - Deep hierarchical structure
  - Repetitive (cyclic) structure
- Not for
  - Systems with a lot of errors and uncertainty
  - Activities with shallow structure

# Basics

---

## Context-Free Grammar

$$G = \langle S, T, N, P \rangle$$

Generic Language	Natural Languages
Start Symbol (S)	Sentences
Set of Terminal Symbols (T)	Words
Set of Non-Terminal Symbols (N)	Parts of Speech
Set of Production Rules (P)	Syntax Rules

# Parsing with a grammar

$S \rightarrow NP VP$	(0.8)	$PP \rightarrow PREP NP$	(1.0)
$S \rightarrow VP$	(0.2)	$PREP \rightarrow like$	(1.0)
$NP \rightarrow NOUN$	(0.4)	$VERB \rightarrow swat$	(0.2)
$NP \rightarrow NOUN PP$	(0.4)	$VERB \rightarrow flies$	(0.4)
$NP \rightarrow NOUN NP$	(0.2)	$VERB \rightarrow like$	(0.4)
$VP \rightarrow VERB$	(0.3)	$NOUN \rightarrow swat$	(0.05)
$VP \rightarrow VERB NP$	(0.3)	$NOUN \rightarrow flies$	(0.45)
$VP \rightarrow VERB PP$	(0.2)	$NOUN \rightarrow ants$	(0.5)
$VP \rightarrow VERB NP PP$	(0.2)		

*swat*

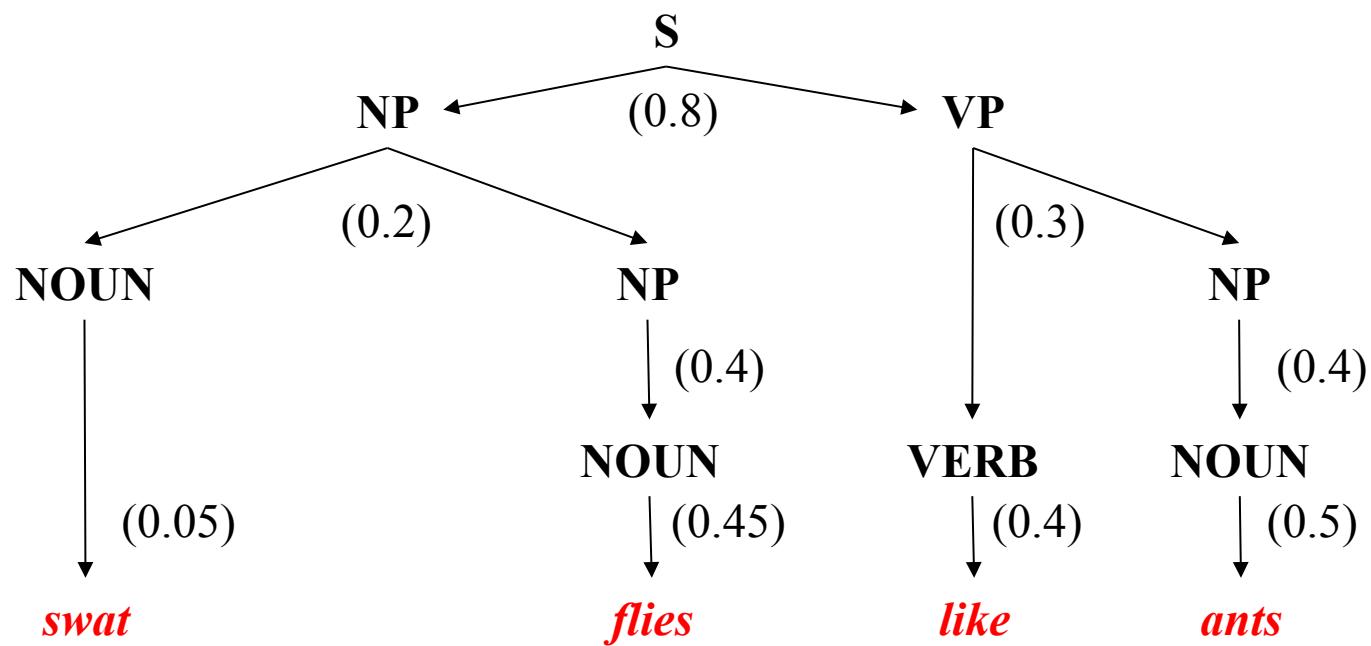
*flies*

*like*

*ants*

# Parsing with a grammar

$S \rightarrow NP VP$	(0.8)	$PP \rightarrow PREP NP$	(1.0)
$S \rightarrow VP$	(0.2)	$PREP \rightarrow like$	(1.0)
$NP \rightarrow NOUN$	(0.4)	$VERB \rightarrow swat$	(0.2)
$NP \rightarrow NOUN PP$	(0.4)	$VERB \rightarrow flies$	(0.4)
$NP \rightarrow NOUN NP$	(0.2)	$VERB \rightarrow like$	(0.4)
$VP \rightarrow VERB$	(0.3)	$NOUN \rightarrow swat$	(0.05)
$VP \rightarrow VERB NP$	(0.3)	$NOUN \rightarrow flies$	(0.45)
$VP \rightarrow VERB PP$	(0.2)	$NOUN \rightarrow ants$	(0.5)
$VP \rightarrow VERB NP PP$	(0.2)		



# Video analysis with CFGs

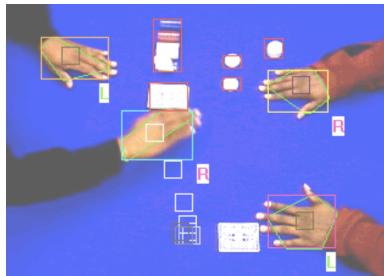
---



The “Inverse Hollywood problem”:  
From video to scripts and storyboards via causal analysis.  
**Brand 1997**

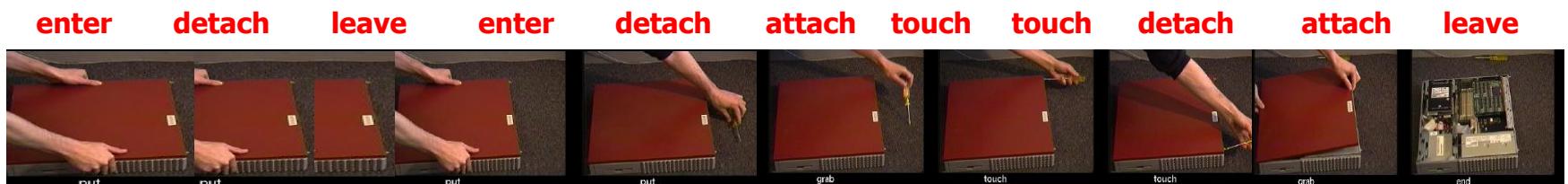


Action Recognition using Probabilistic Parsing.  
**Bobick and Ivanov 1998**



Recognizing Multitasked Activities from Video using  
Stochastic Context-Free Grammar.  
**Moore and Essa 2001**

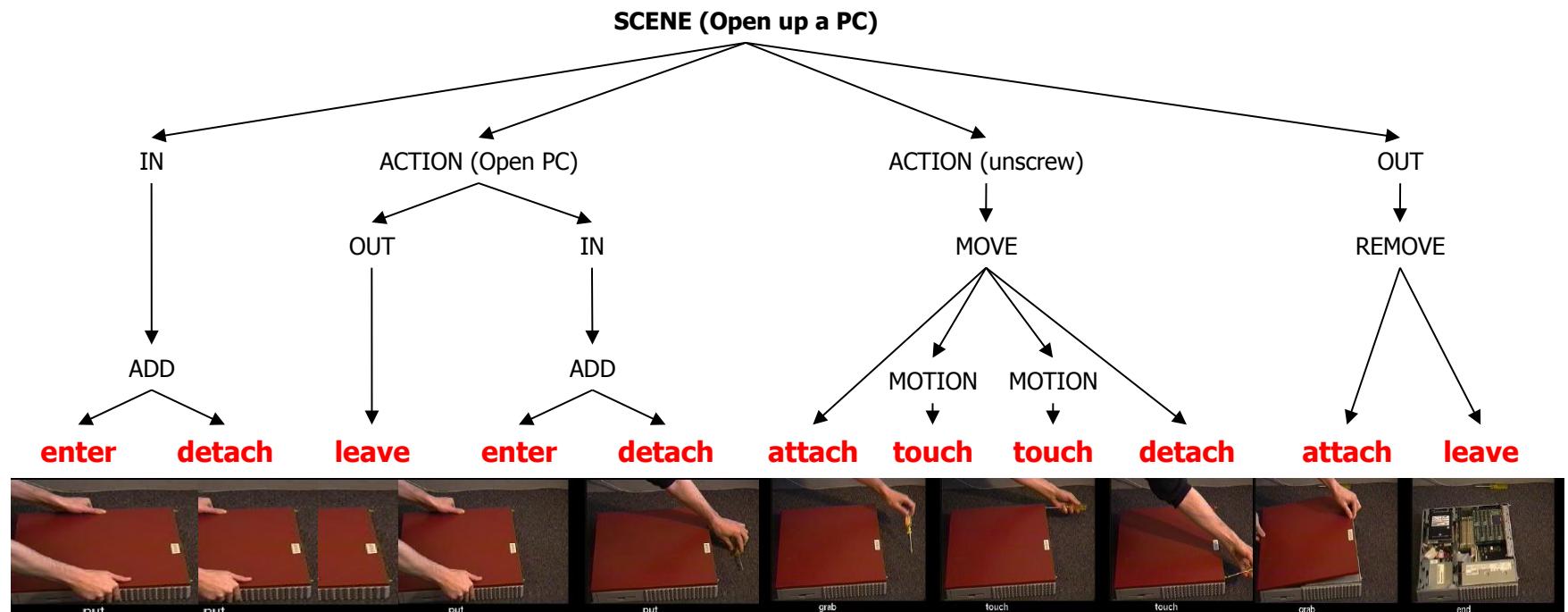
# CFG for human activities



scene	→	in action* out
action	→	motion   move   {out in}
in	→	ENTER   add
out	→	LEAVE   remove
add	→	ENTER motion* DETACH
remove	→	ATTACH motion* LEAVE
move	→	ATTACH motion+ DETACH
motion	→	SHIFT   TOUCH   BUMP

M. Brand. The "Inverse Hollywood Problem":  
From video to scripts and storyboards  
via causal analysis. AAAI 1997.

# Parse tree



scene	$\rightarrow$	in action* out
action	$\rightarrow$	motion   move   {out in}
in	$\rightarrow$	ENTER   add
out	$\rightarrow$	LEAVE   remove
add	$\rightarrow$	ENTER motion* DETACH
remove	$\rightarrow$	ATTACH motion* LEAVE
move	$\rightarrow$	ATTACH motion+ DETACH
motion	$\rightarrow$	SHIFT   TOUCH   BUMP

- Deterministic low-level primitive detection
- Deterministic parsing

M. Brand. The "Inverse Hollywood Problem": From video to scripts and storyboards via causal analysis. AAAI 1997.

# Stochastic CFGs

---

Action Recognition using Probabilistic Parsing.  
Bobick and Ivanov 1998



$G_{square} :$			
SQUARE	$\rightarrow$	RH	[0.5]
		LH	[0.5]
RH	$\rightarrow$	TOP UD BOT DU	[1.0]
LH	$\rightarrow$	BOT DU TOP UD	[1.0]
TOP	$\rightarrow$	LR	[0.5]
		RL	[0.5]
BOT	$\rightarrow$	RL	[0.5]
		LR	[0.5]
LR	$\rightarrow$	left-right	[1.0]
UD	$\rightarrow$	up-down	[1.0]
RL	$\rightarrow$	right-left	[1.0]
DU	$\rightarrow$	down-up	[1.0]

# Gesture analysis with CFGs

---

Primitive recognition with HMMs



---

left-right



---

up-down



---

right-left



---

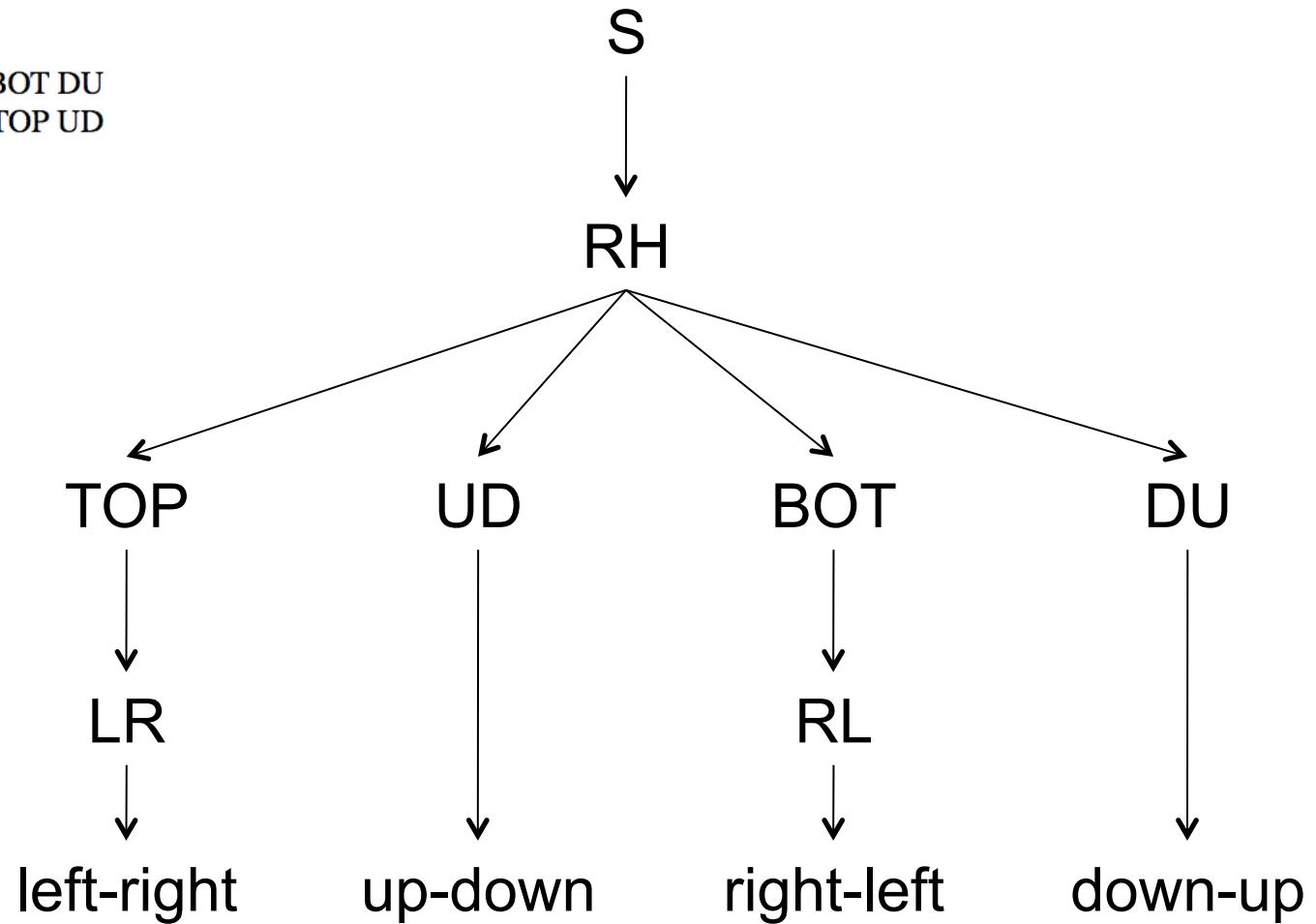
down-up



# Parse Tree

$G_{square} :$

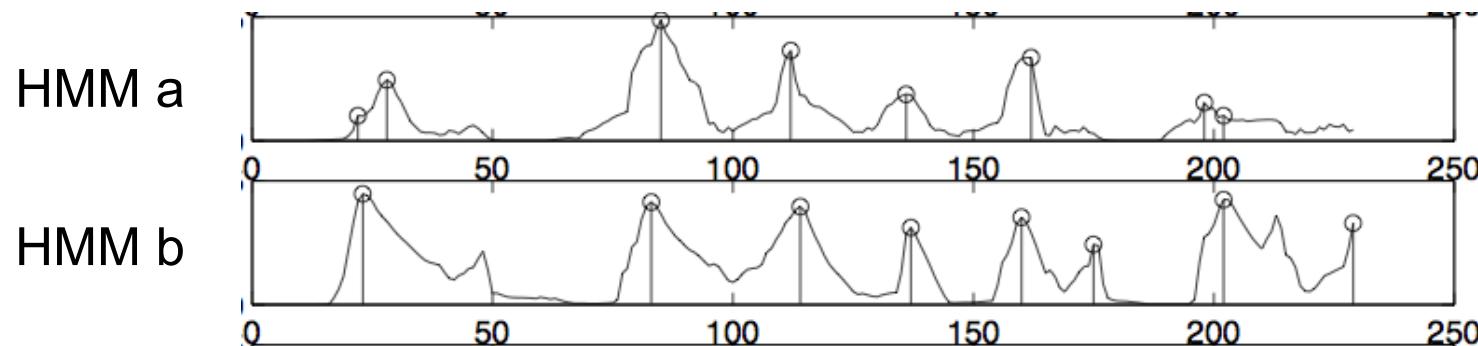
SQUARE	$\rightarrow$	RH
		LH
RH	$\rightarrow$	TOP UD BOT DU
LH	$\rightarrow$	BOT DU TOP UD
TOP	$\rightarrow$	LR
		RL
BOT	$\rightarrow$	RL
		LR
LR	$\rightarrow$	left-right
UD	$\rightarrow$	up-down
RL	$\rightarrow$	right-left
DU	$\rightarrow$	down-up



# Errors

---

Likelihood value over time (not discrete symbols)



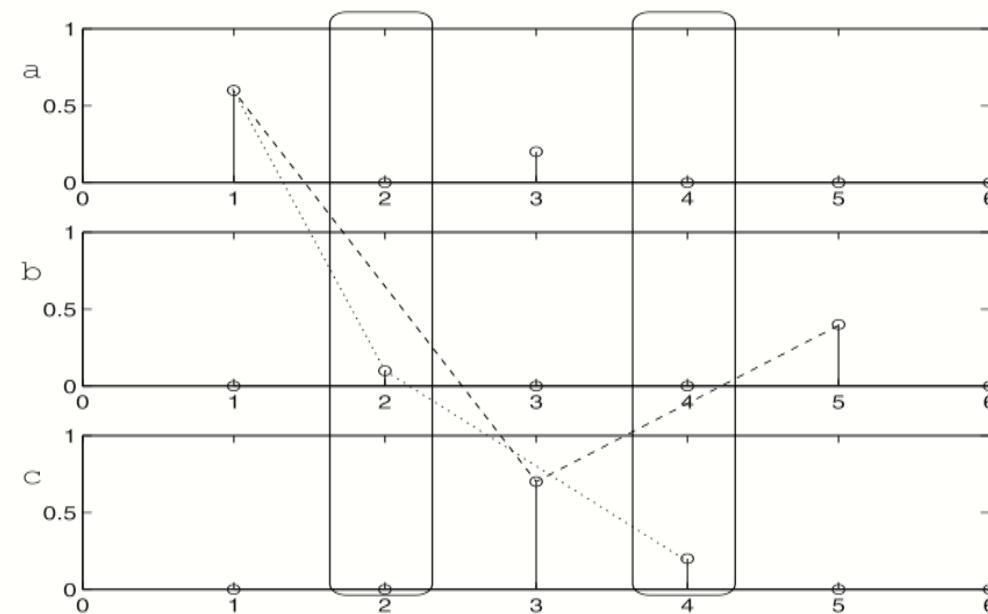
Errors are inevitable...

but the grammar acts as a top-down constraint

# Dealing with uncertainty & errors

- Stolcke-Early (probabilistic) parser
- SKIP rules to deal with **insertion** errors

HMM a



HMM b

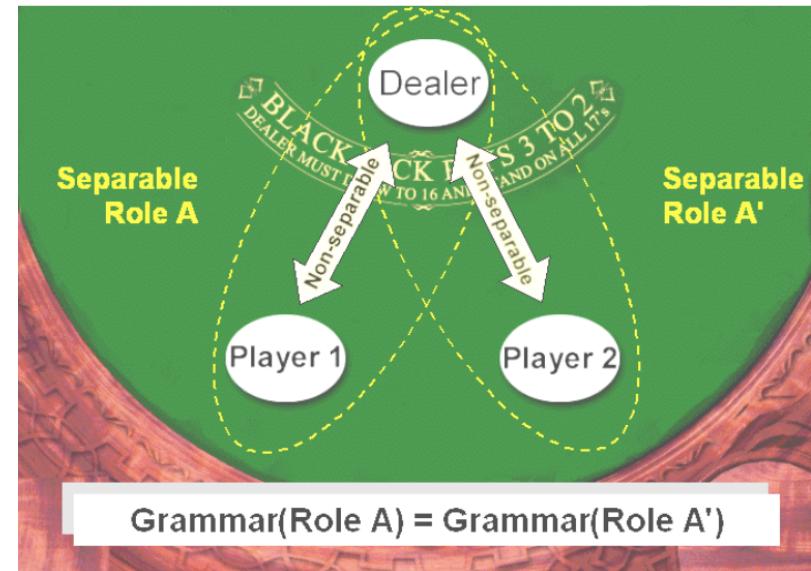
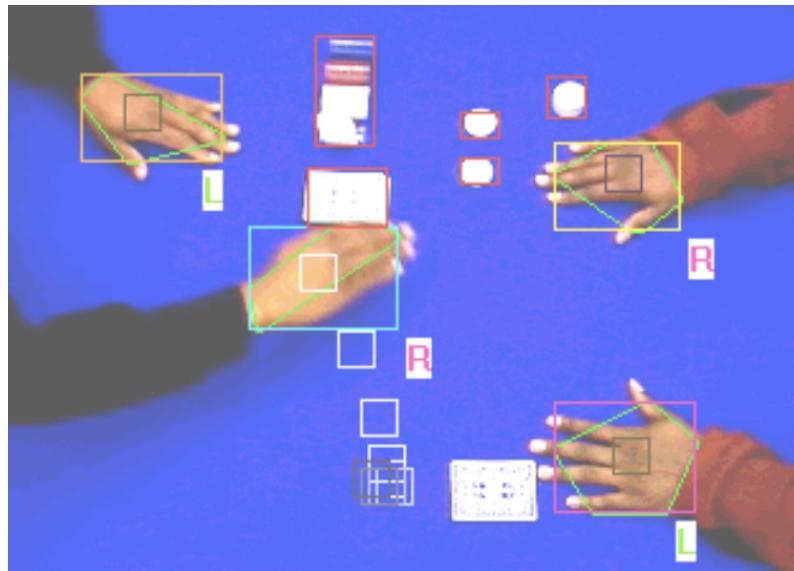


HMM c



# SCFG for Blackjack

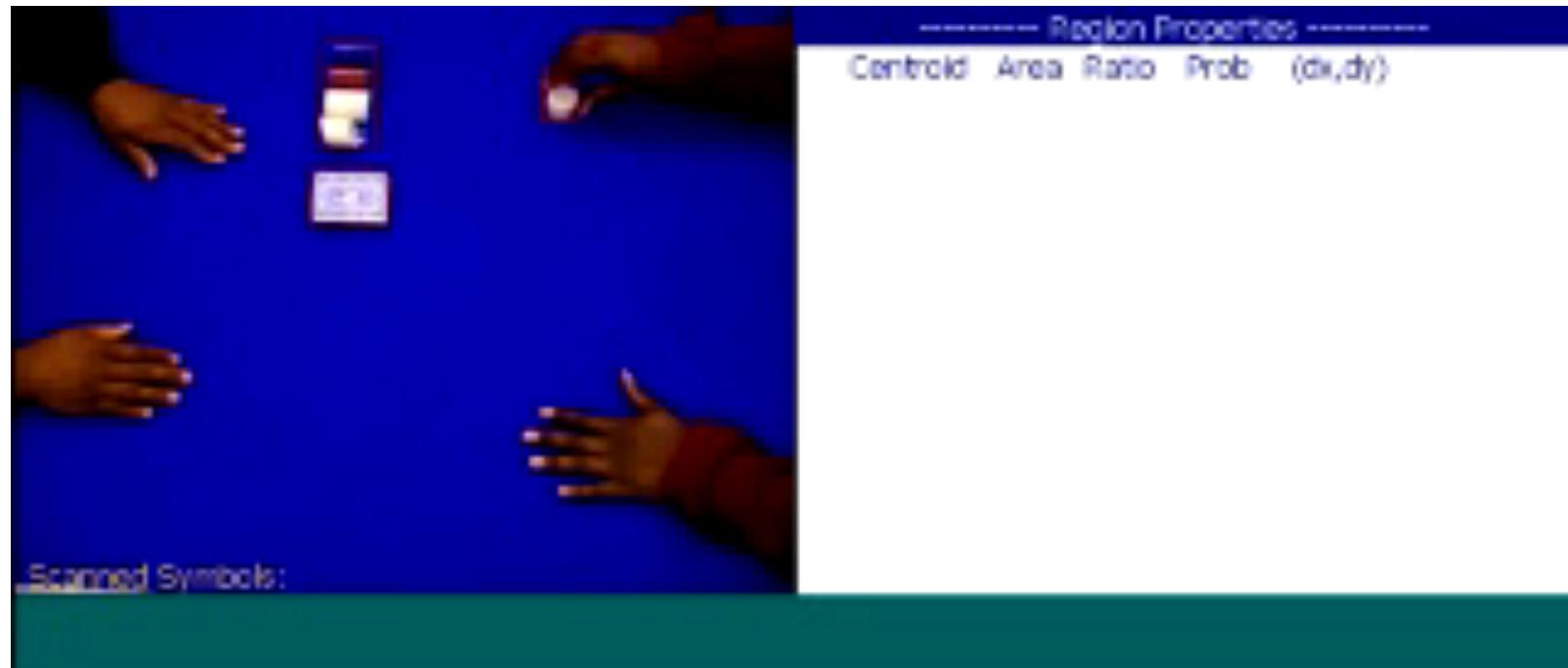
Recognizing Multitasked Activities from Video using  
Stochastic Context-Free Grammar.  
Moore and Essa 2001



- Deals with more complex activities
- Deals with more error types

# extracting primitive actions

---



# Game grammar

Production Rules			Description
<i>S</i>	$\rightarrow AB$	[1.0]	Blackjack $\rightarrow$ "play game" "determine winner"
<i>A</i>	$\rightarrow CD$	[1.0]	play game $\rightarrow$ "setup game" "implement strategy"
<i>B</i>	$\rightarrow EF$	[1.0]	determine winner $\rightarrow$ "eval. strategy" "cleanup"
<i>C</i>	$\rightarrow HI$	[1.0]	setup game $\rightarrow$ "place bets" "deal card pairs"
<i>D</i>	$\rightarrow GK$	[1.0]	implement strategy $\rightarrow$ "player strategy"
<i>E</i>	$\rightarrow LKM$	[0.6]	eval. strategy $\rightarrow$ "dealer down-card" "dealer hits" "player down-card"
	$\rightarrow LM$	[0.4]	eval. strategy $\rightarrow$ "dealer down-card" "player down-card"
<i>F</i>	$\rightarrow NO$	[0.5]	cleanup $\rightarrow$ "settle bet" "recover card"
	$\rightarrow ON$	[0.5]	$\rightarrow$ "recover card" "settle bet"
<i>G</i>	$\rightarrow J$	[0.8]	player strategy $\rightarrow$ "Basic Strategy"
	$\rightarrow Hf$	[0.1]	$\rightarrow$ "Splitting Pair"
	$\rightarrow bfffH$	[0.1]	$\rightarrow$ "Doubling Down"
<i>H</i>	$\rightarrow l$	[0.5]	place bets
	$\rightarrow lH$	[0.5]	
<i>I</i>	$\rightarrow ffI$	[0.5]	deal card pairs
	$\rightarrow ee$	[0.5]	
<i>J</i>	$\rightarrow f$	[0.8]	Basic strategy
	$\rightarrow fJ$	[0.2]	
<i>K</i>	$\rightarrow e$	[0.6]	house hits
	$\rightarrow eK$	[0.4]	
<i>L</i>	$\rightarrow ae$	[1.0]	Dealer downcard
<i>M</i>	$\rightarrow dh$	[1.0]	Player downcard
<i>N</i>	$\rightarrow k$	[0.16]	settle bet
	$\rightarrow kN$	[0.16]	
	$\rightarrow j$	[0.16]	
	$\rightarrow jN$	[0.16]	
	$\rightarrow i$	[0.18]	
	$\rightarrow iN$	[0.18]	
<i>O</i>	$\rightarrow a$	[0.25]	recover card
	$\rightarrow aO$	[0.25]	
	$\rightarrow b$	[0.25]	
	$\rightarrow bO$	[0.25]	

Symbol	Domain-Specific Events (Terminals)
<i>a</i>	dealer removed card from house
<i>b</i>	dealer removed card from player
<i>c</i>	player removed card from house
<i>d</i>	player removed card from player
<i>e</i>	dealer added card to house
<i>f</i>	dealer dealt card to player
<i>g</i>	player added card to house
<i>h</i>	player added card to player
<i>i</i>	dealer removed chip
<i>j</i>	player removed chip
<i>k</i>	dealer pays player chip
<i>l</i>	player bets chip

# Dealing with errors

---

- Ungrammatical strings cause parser to fail
- Account for errors with multiple hypothesis
  - Insertion, deletion, substitution
- Issues
  - How many errors should we tolerate?
  - Potentially exponential hypothesis space
  - Ungrammatical strings: vision problem or illegal activity?

# Observations

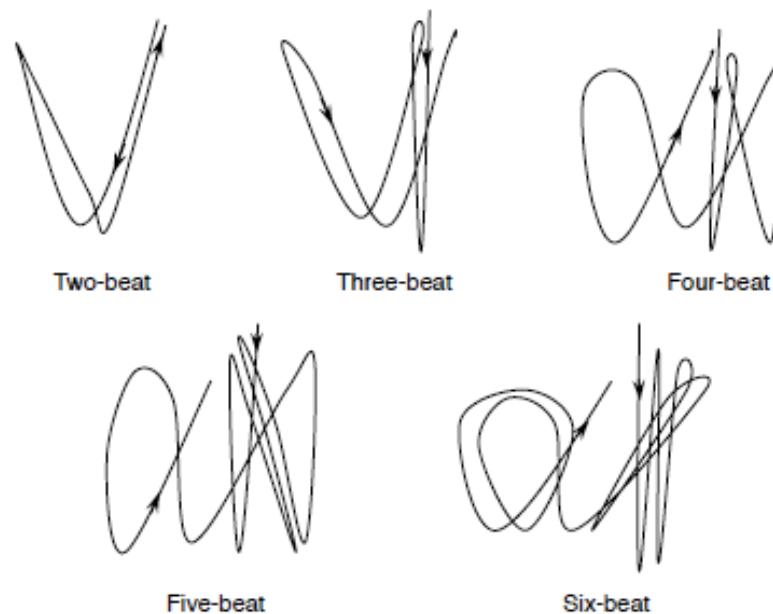
---

- CFGs good for structured activities
  - Can incorporate uncertainty in observations
  - Natural contextual prior for recognizing errors
- 
- Not clear how to deal with errors
  - Assumes ‘good’ action classifiers
  - Need to define grammar manually

**Can we learn the grammar from data?**

# Heuristic Grammatical Induction

---



1. Lexicon learning
  - Learn HMMs
  - Cluster HMMs
2. Convert video to string
3. Learn Grammar

Unsupervised Analysis of Human Gestures. Wang et al 2001

# COMPRESSIVE

a b c d a b c d b c d a b a b

$$\arg \max_{\lambda} \Delta DL = \arg \max_{\lambda} \{ M \times N - (M + 1) - N \}$$

length   occurrence   new rule   new symbol  
 deletion of   insertion of  
 substring   new rule

substring	M	N	$\Delta DL$
<b>ab</b>	2	4	1
<b>cd</b>	2	3	0
<b>bcd</b>	3	3	2
<b>abcd</b>	3	2	1

# On-Line and Off-Line Heuristics for Inferring Hierarchies of Repetitions in Sequences. Nevill-Manning 2000.

# example

---

$S \rightarrow a b c d a b c d b c d a b a b$   
(DL=16)

$A \rightarrow b c d$

$S \rightarrow a A a A A a b a b$   
(DL=14)

Repeat until compression becomes 0.

# Critical assumption

---

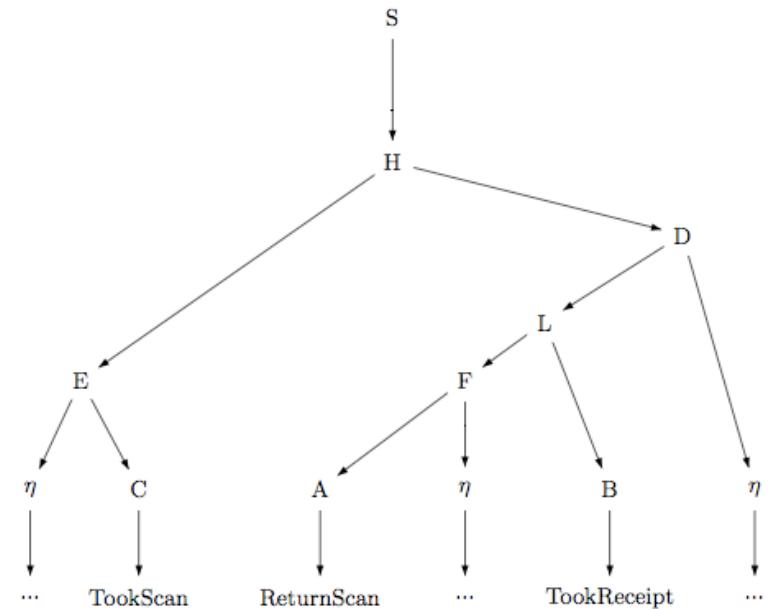
- No uncertainty
- No errors
  - insertions
  - deletions
  - substitution

**Can we learn grammars despite errors?**

# Learning with noise

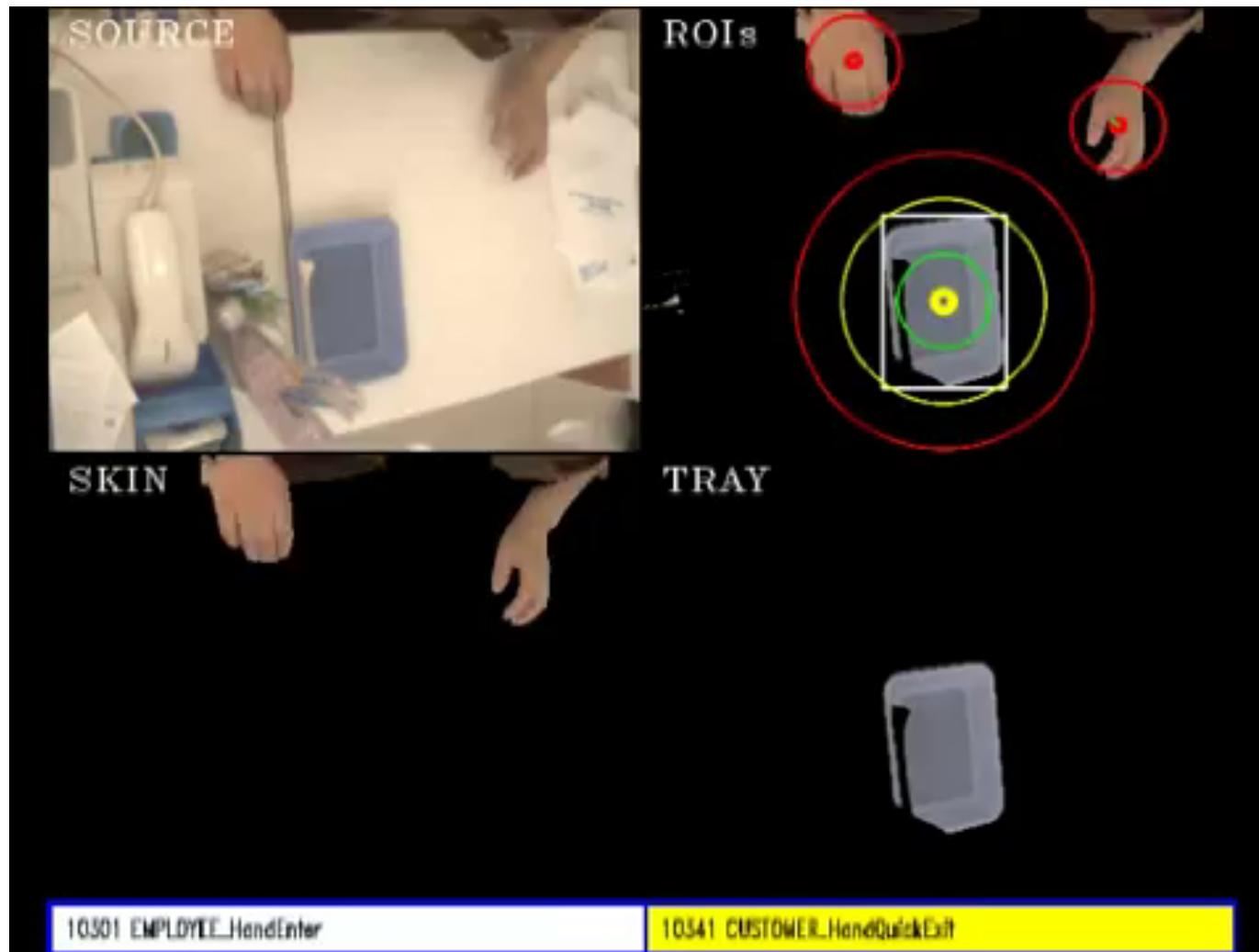
---

Can we learn the basic structure of a transaction?



Recovering the basic structure of human activities from  
noisy video-based symbol strings. Kitani et al 2008.

# extracting primitives



Recovering the basic structure of human activities from noisy video-based symbol strings. Kitani et al 2008.

# Underlying structure?

---

D → a x b y c a b x c y a b c x

# Underlying structure?

---

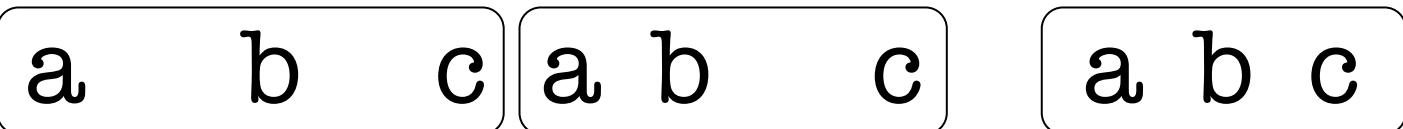
$D \rightarrow a \ x \ b \ y \ c \ a \ b \ x \ c \ y \ a \ b \ c \ x$

$D \rightarrow a \quad b \quad c \ a \ b \quad c \quad a \ b \ c$

# Underlying structure?

---

$D \rightarrow a \ x \ b \ y \ c \ a \ b \ x \ c \ y \ a \ b \ c \ x$

$D \rightarrow$   a      b      c      a    b      c      a    b    c

# Underlying structure?

---

$D \rightarrow a \ x \ b \ y \ c \ a \ b \ x \ c \ y \ a \ b \ c \ x$

$D \rightarrow a \quad b \quad c \ a \ b \quad c \quad a \ b \ c$

$A \rightarrow a \ b \ c$

Simple grammar

$D \rightarrow A \ A \ A$

Efficient compression

# Information Theory Problem (MDL)

---

$$\hat{G} = \arg \min_G \{ DL(G) + \underset{\text{Model complexity}}{DL(D|G)} \}$$

# Information Theory Problem (MDL)

---

$$\hat{G} = \arg \min_G \{ DL(G) + \text{Model complexity} \cdot DL(D|G) \}$$

$$\begin{aligned} DL(G) &= -\log p(G) \\ \text{Model complexity} &= -\log p(\theta_S, G_S) \\ &= -\log p(\theta_S|G_S) - \log p(G_S) \\ &= DL(\theta_S|G_S) - \text{Grammar parameters} \cdot DL(G_S) \end{aligned} \quad \begin{aligned} & & \text{Grammar structure} \end{aligned}$$

# Information Theory Problem (MDL)

---

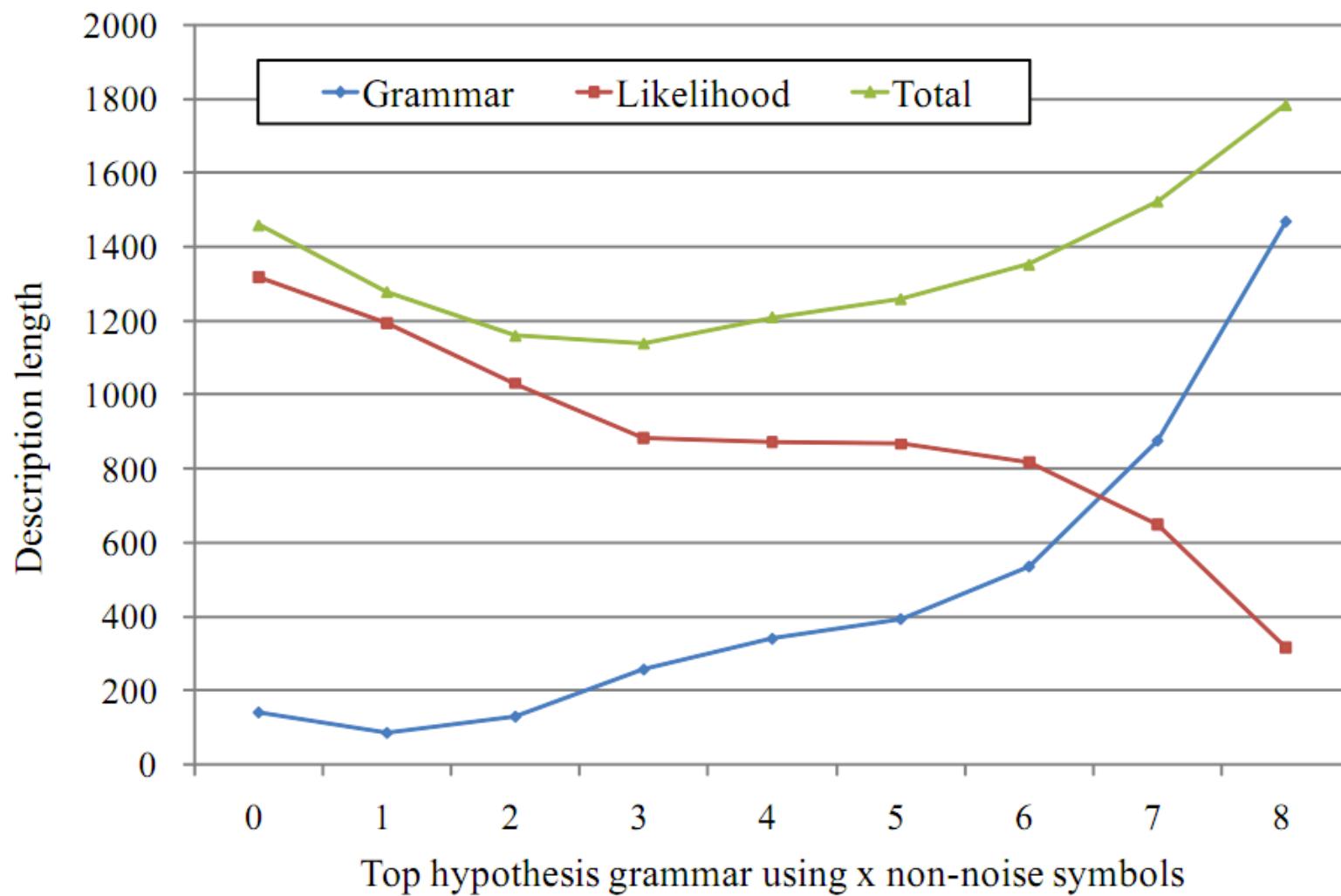
$$\hat{G} = \arg \min_G \{ DL(G) + \underset{\text{Model complexity}}{DL(D|G)} \}$$

$$\begin{aligned} \underset{\text{Model complexity}}{DL(G)} &= -\log p(G) \\ &= -\log p(\theta_S, G_S) \\ &= -\log p(\theta_S|G_S) - \log p(G_S) \\ &= \underset{\text{Grammar parameters}}{DL(\theta_S|G_S)} - \underset{\text{Grammar structure}}{DL(G_S)} \end{aligned}$$

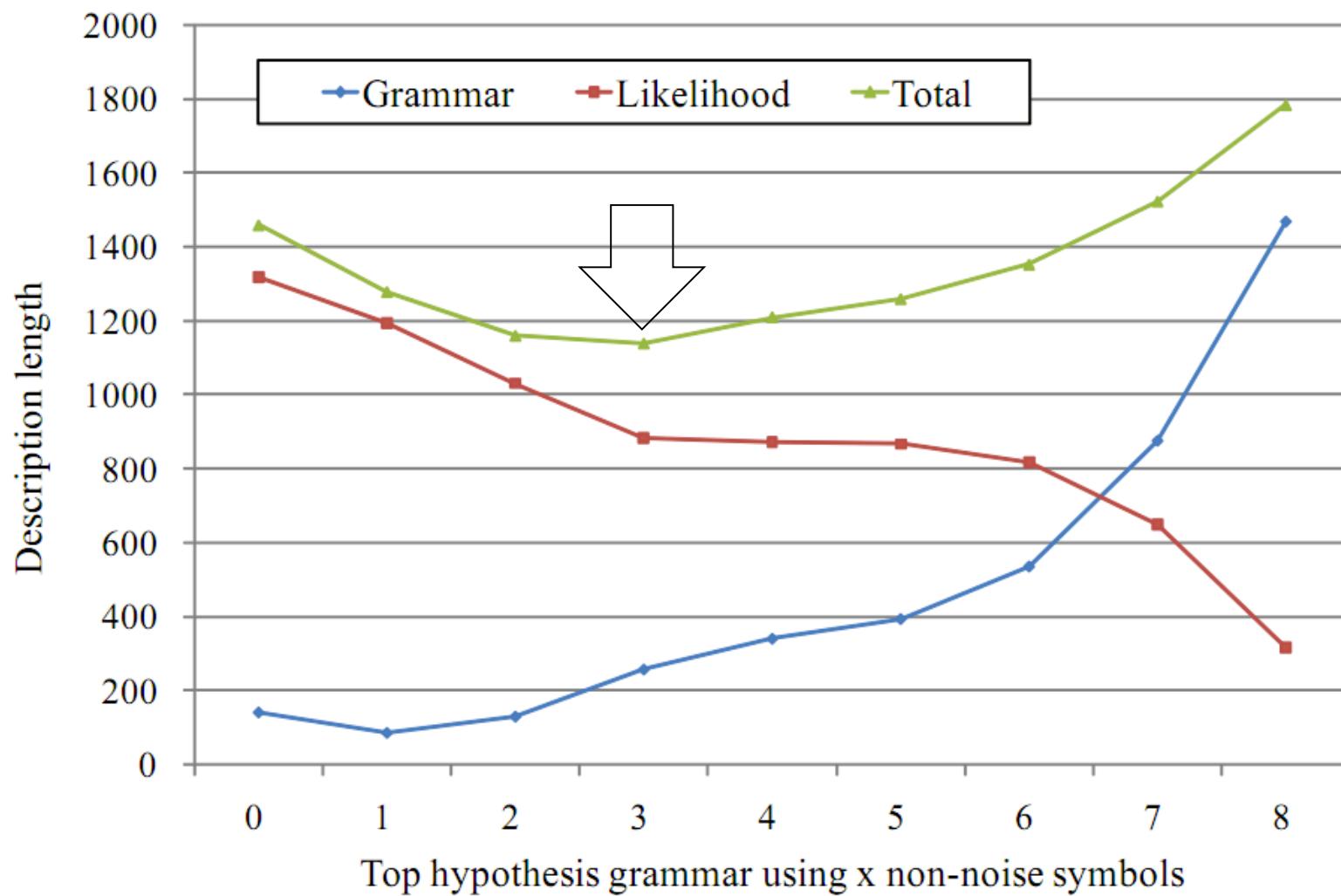
$$\underset{\text{Data compression}}{DL(D|G)} = -\log p(D|G)$$

Likelihood  
(inside probabilities)

# Minimum Description Length

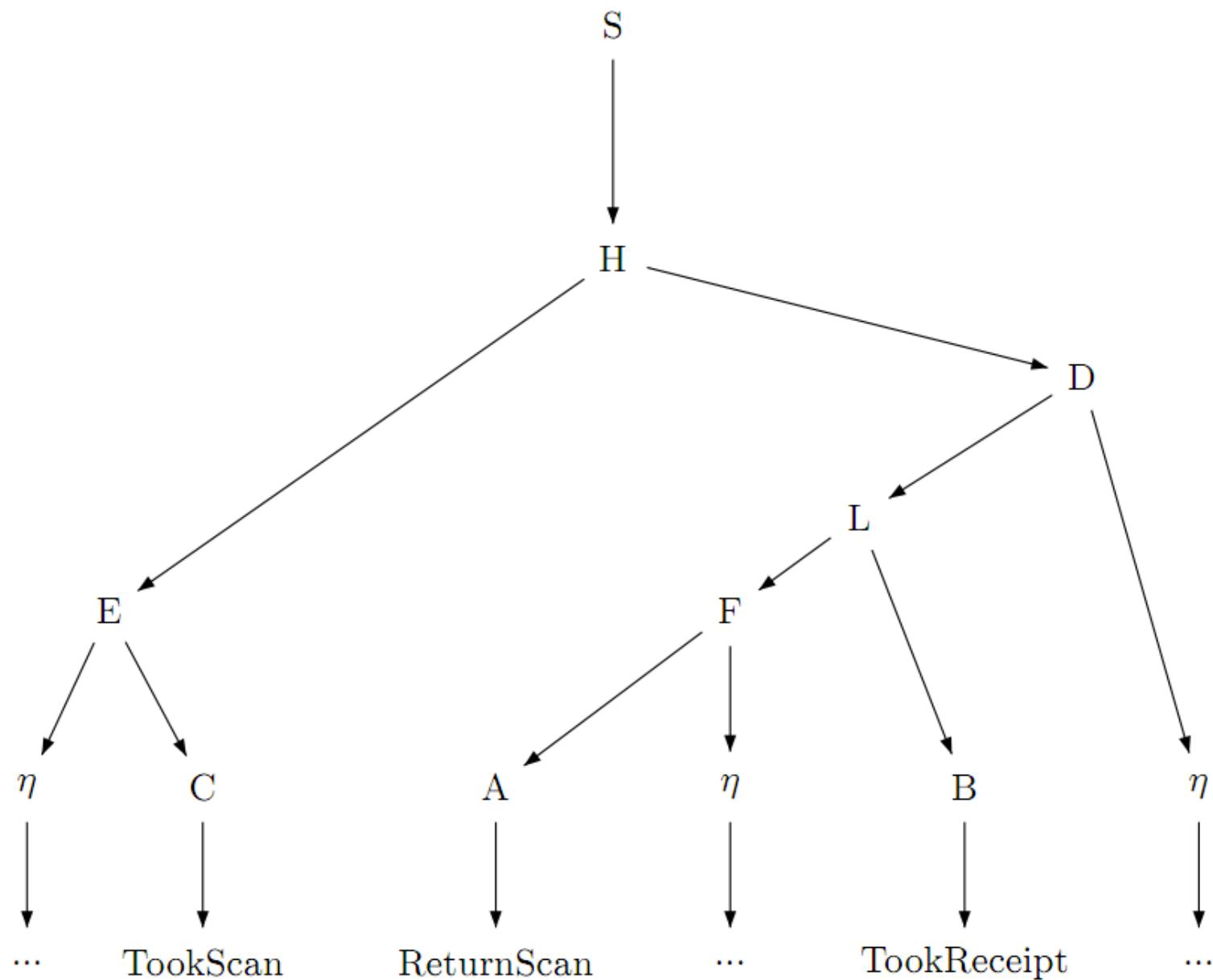


# Minimum Description Length



Recovering the basic structure of human activities from noisy video-based symbol strings. Kitani et al 2008.

$S \rightarrow D$		(0.02)	$D \rightarrow L$	$\eta$	(1.000)
$S \rightarrow H$		(0.16)	$E \rightarrow \eta$	C	(1.000)
$S \rightarrow G$		(0.18)	$F \rightarrow A$	$\eta$	(1.000)
$S \rightarrow N$	$\eta$	(0.04)	$G \rightarrow C$	D	(1.000)
$S \rightarrow J$		(0.13)	$H \rightarrow E$	D	(1.000)
$S \rightarrow Q$		(0.05)	$I \rightarrow *$	B	$\eta$ (1.000)
$S \rightarrow \eta$		(0.02)	$J \rightarrow C$	F	(1.000)
$S \rightarrow N$		(0.02)	$K \rightarrow *$	D	(1.000)
$S \rightarrow R$		(0.05)	$L \rightarrow F$	B	(1.000)
$S \rightarrow J$	B	(0.02)	$M \rightarrow C$	*	(1.000)
$S \rightarrow M$	L	(0.04)	$N \rightarrow E$	A	B (1.000)
$S \rightarrow M$	A H	(0.02)	$O \rightarrow E$	*	(1.000)
$S \rightarrow C$	K	(0.04)	$P \rightarrow E$	I	(1.000)
$S \rightarrow C$	A M	F	$Q \rightarrow E$	K	(1.000)
$S \rightarrow O$	F	(0.02)	$R \rightarrow E$	L	(1.000)
$S \rightarrow M$		(0.02)	$\eta \rightarrow \eta$	$\eta$	(0.309)
$S \rightarrow O$	L	(0.02)	$\eta \rightarrow$	CUS_AddMoney	(0.153)
$S \rightarrow P$		(0.05)	$\eta \rightarrow$	CUS_MovedTray	(0.006)
$S \rightarrow I$		(0.04)	$\eta \rightarrow$	CUS_RemMoney	(0.003)
$S \rightarrow K$		(0.04)	$\eta \rightarrow$	EMP_HandReturn	(0.080)
$A \rightarrow$	EMP_ReturnedScanner	(1.00)	$\eta \rightarrow$	EMP_Interaction	(0.275)
$B \rightarrow$	EMP_TookReceipt	(1.00)	$\eta \rightarrow$	EMP_MovedTray	(0.028)
$C \rightarrow$	EMP_TookScanner	(1.00)	$\eta \rightarrow$	EMP_RemMoney	(0.147)



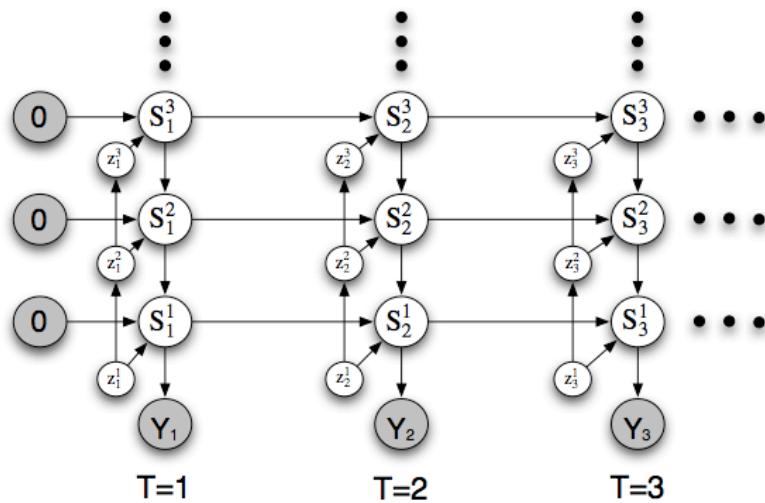
Recovering the basic structure of human activities from noisy video-based symbol strings. Kitani et al 2008.

# Conclusions

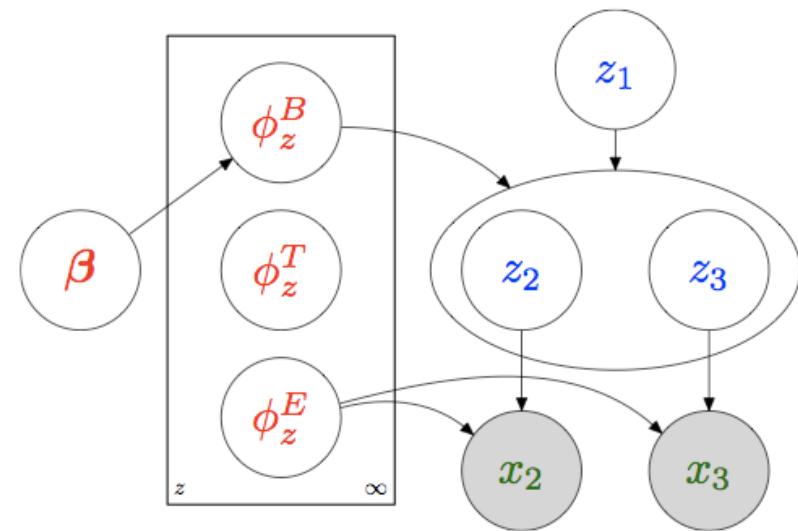
---

- Possible to learn basic structure
- Robust to errors  
(insertion, deletion, substitution)
- Need a lot of training data
- Computational complexity

# Bayesian Approaches



Infinite Hierarchical Hidden Markov Models.  
Heller et al 2009.



The Infinite PCFG using Hierarchical Dirichlet Processes.  
Liang et al 2007.

# Take home message

## Hierarchical Syntactic Models

---

- Useful for activities with:
  - Deep hierarchical structure
  - Repetitive (cyclic) structure
- Not for
  - Systems with a lot of errors and uncertainty
  - Activities with weak structure

---

# **Statistical Approaches**

# Using a hierarchical statistical approach

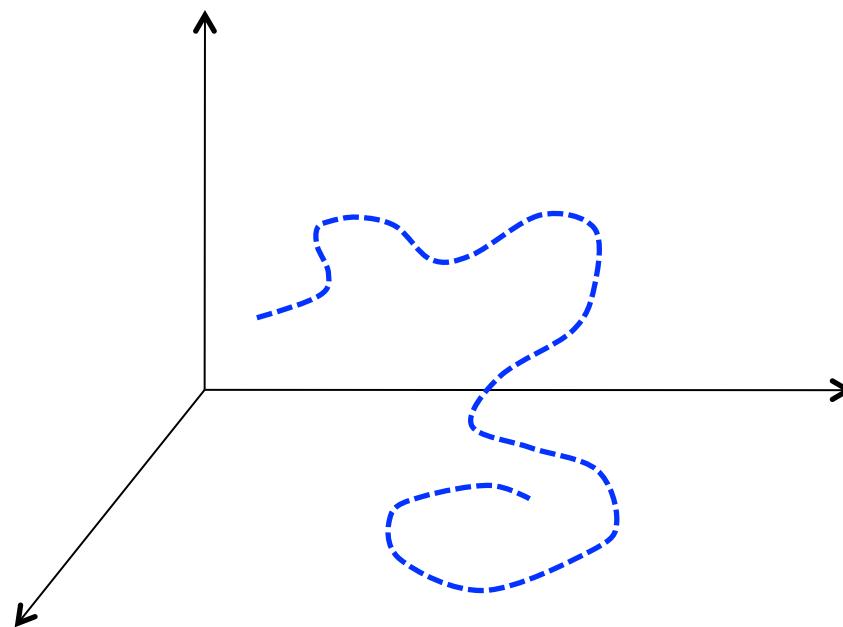
---

- Use when
  - Low-level action detectors are noisy
  - Structure of activity is sequential
  - Integrating dynamics
- Not for
  - Activities with deep hierarchical structure
  - Activities with complex temporal structure

# Statistical (State-based) Model

---

Activities as a stochastic path.



What are the underlying dynamics?

# Characteristics

---

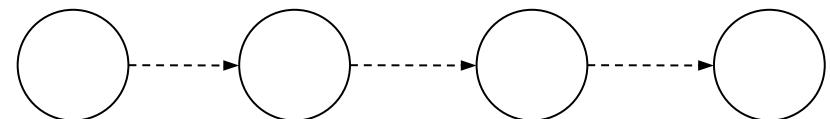
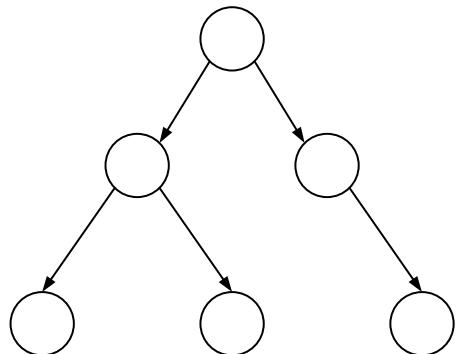
- Strong Markov assumption
  - Strong dynamics prior
  - Robust to uncertainty
- 
- Modifications to account for
    - Hierarchical structure
    - Concurrent structure

# Hierarchical activities

---

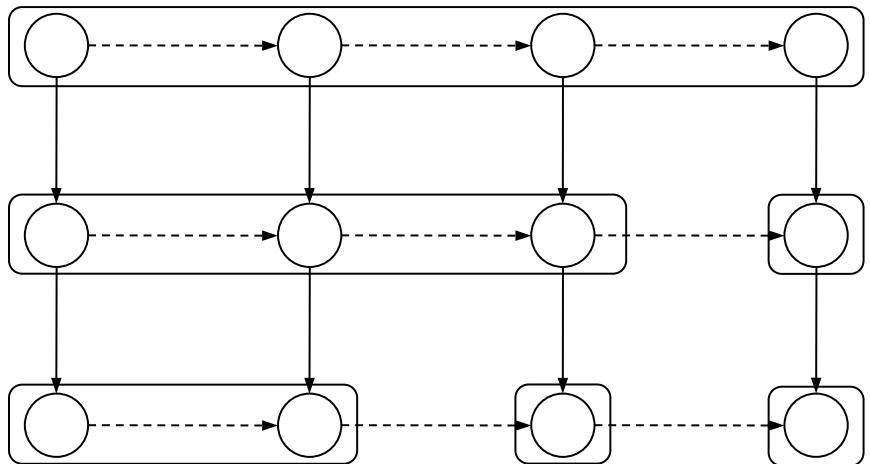
**Problem:**

How do we model  
hierarchical activities?



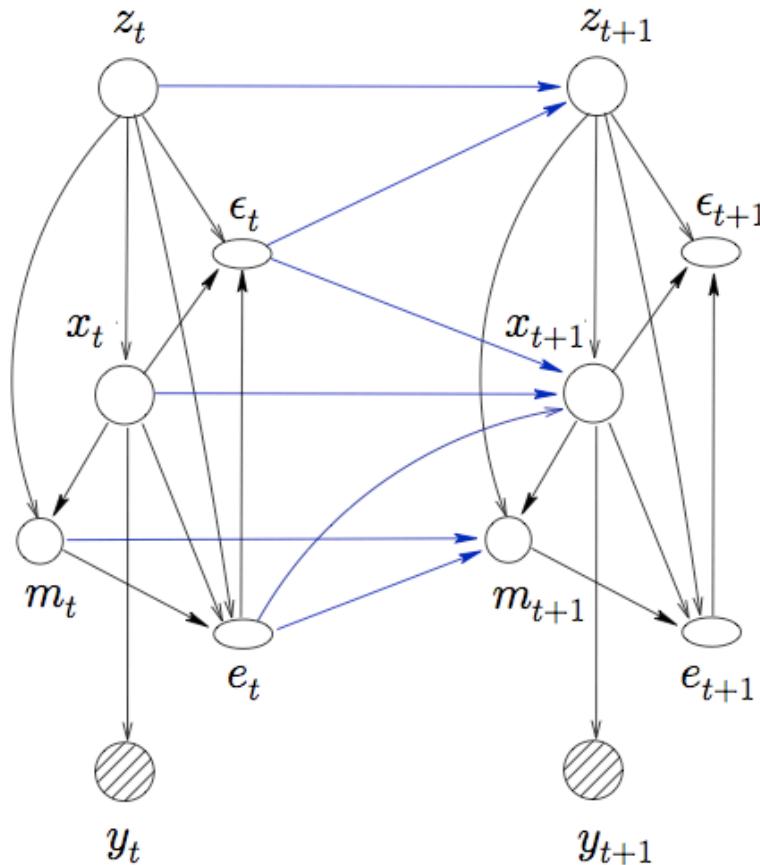
combinatory state space!

**Solution:**  
“**stack**” actions for  
hierarchical activities



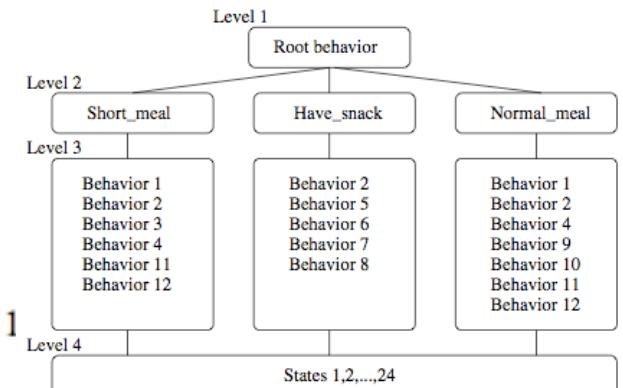
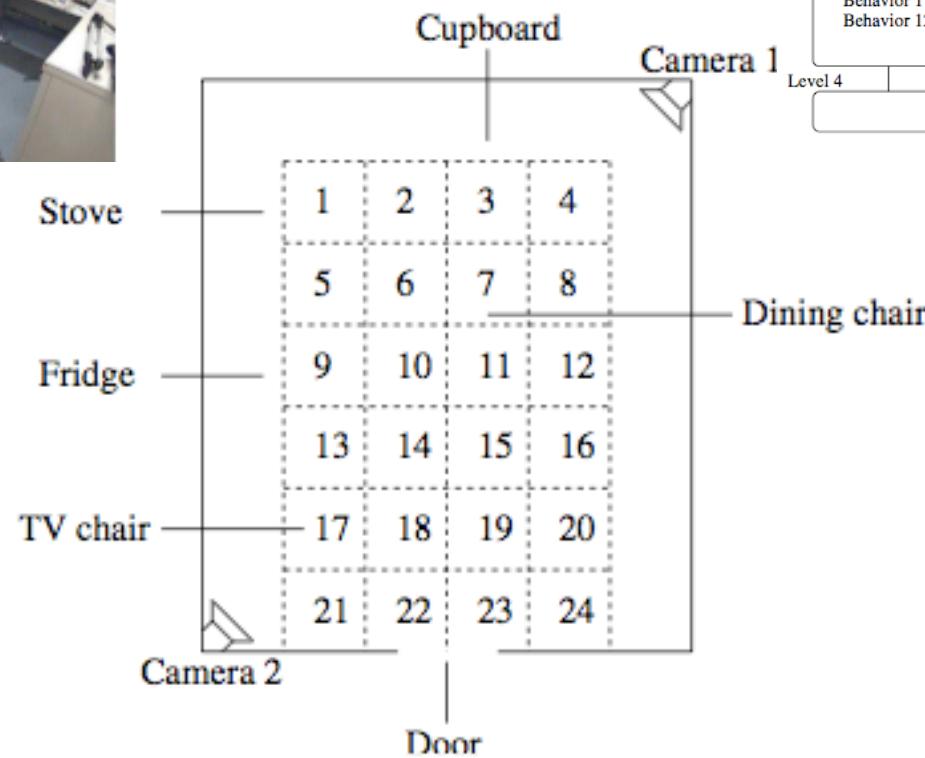
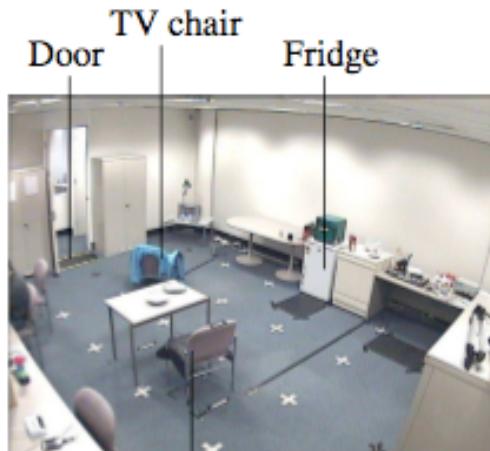
# Hierarchical hidden Markov model

---

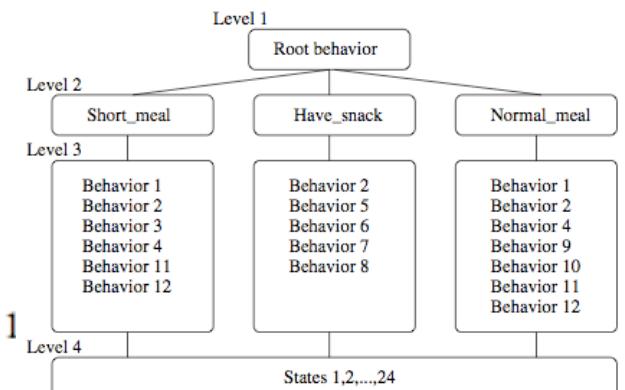
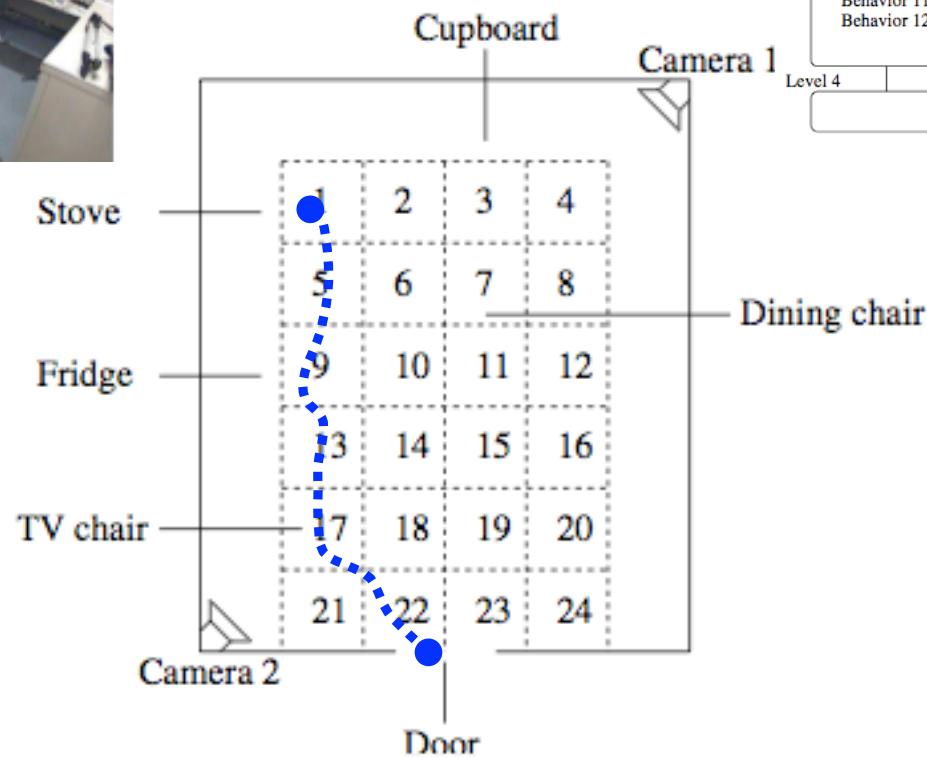
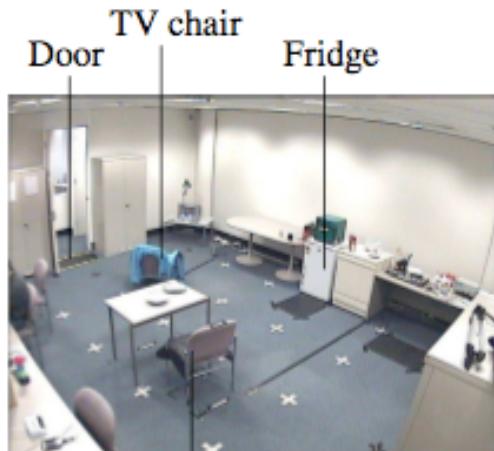


Learning and Detecting Activities from Movement Trajectories Using the  
Hierarchical Hidden Markov Models. Nguyen et al 2005

# Context-free activity grammar



# Context-free activity grammar



# Observations

---

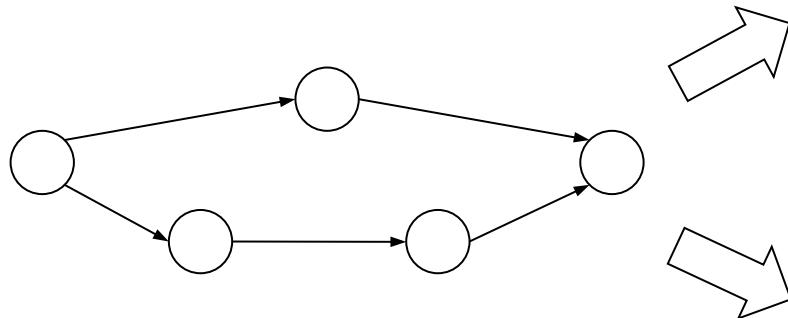
- Tree structures useful for hierarchies
- Tight integration of trajectories with abstract semantic states
- Activities are not always a single sequence  
(ie. they sometimes happen in parallel )

# Concurrent activities

---

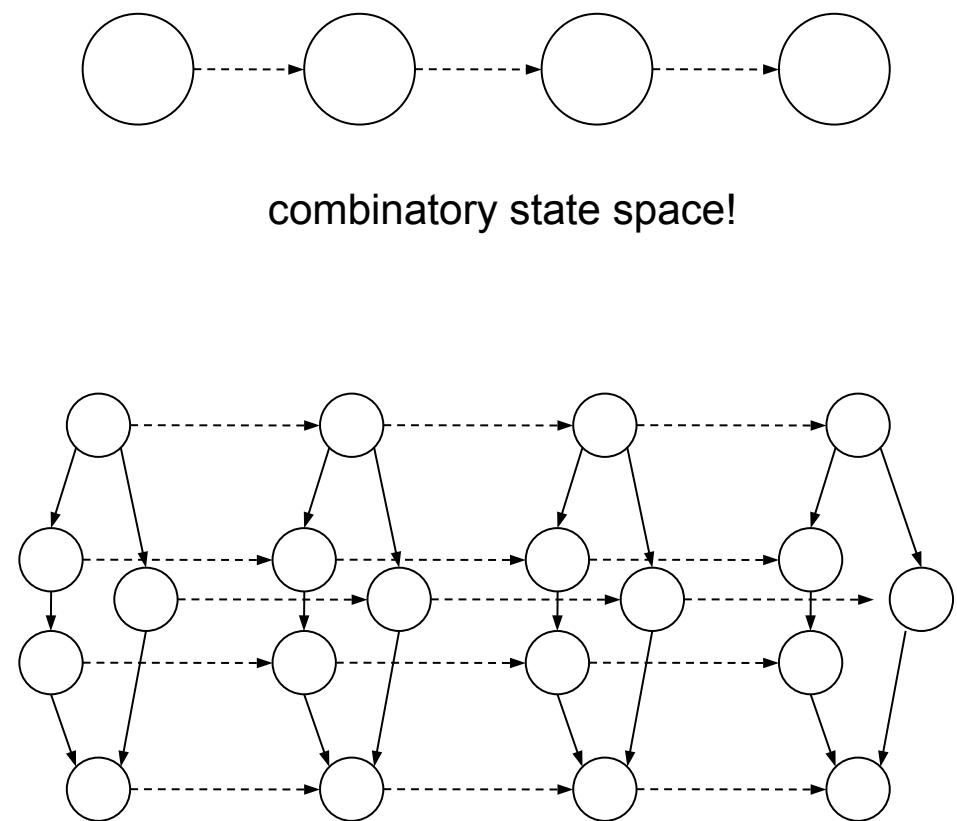
**Problem:**

How do we model  
concurrent activities?

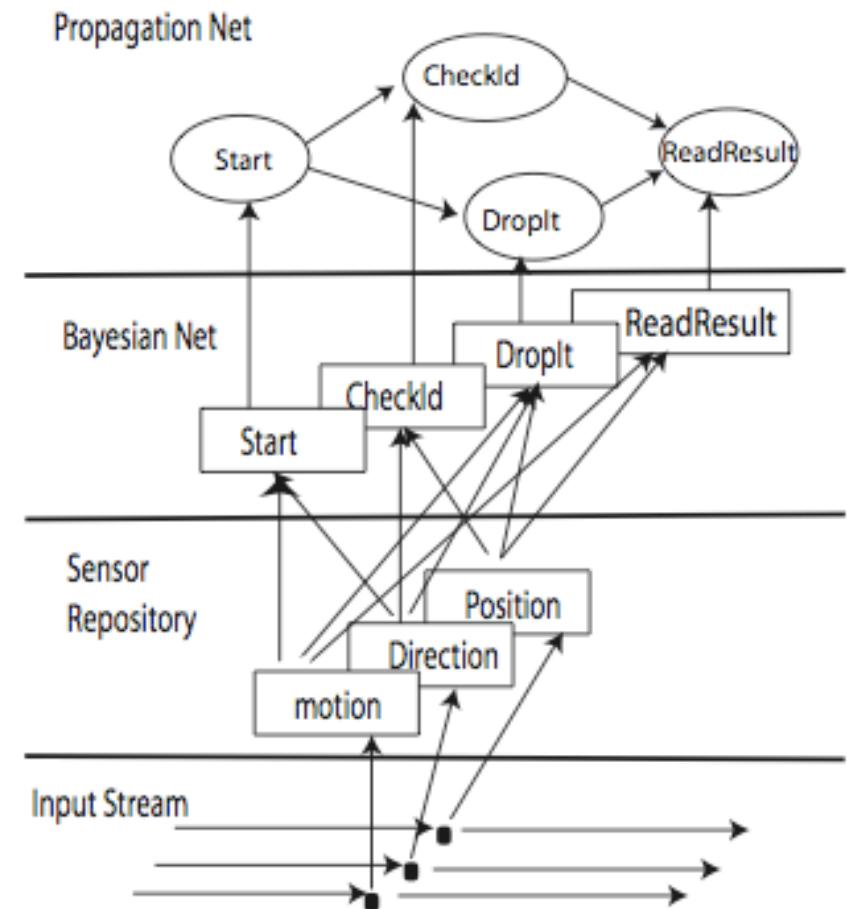
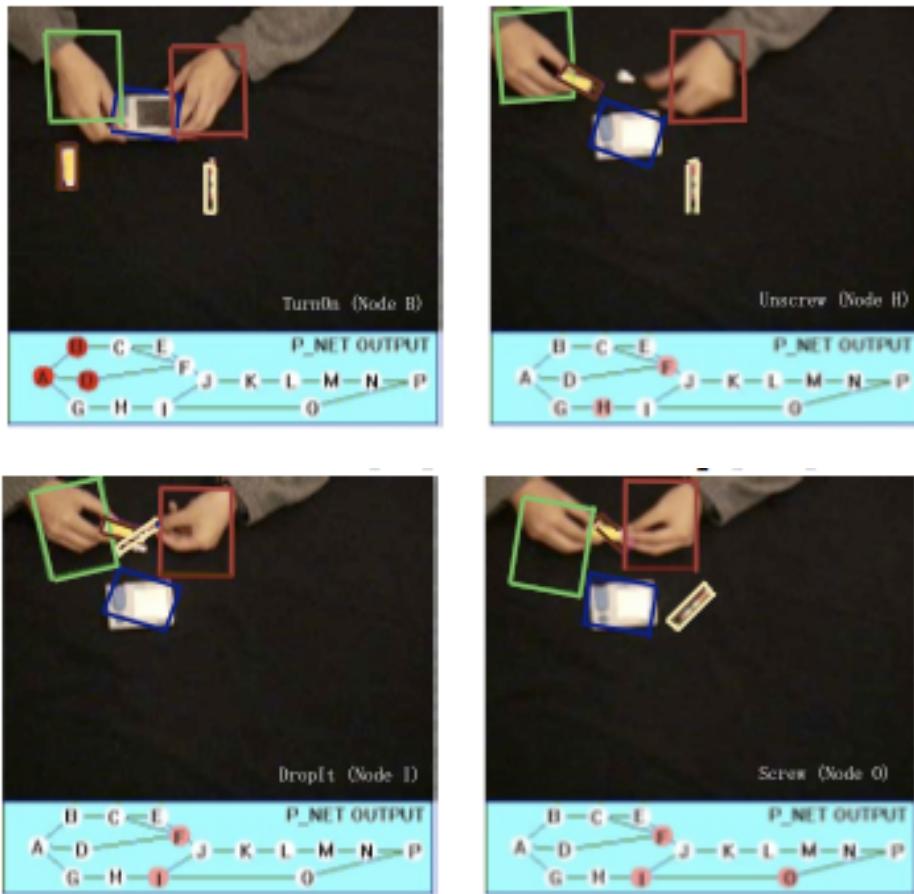


combinatory state space!

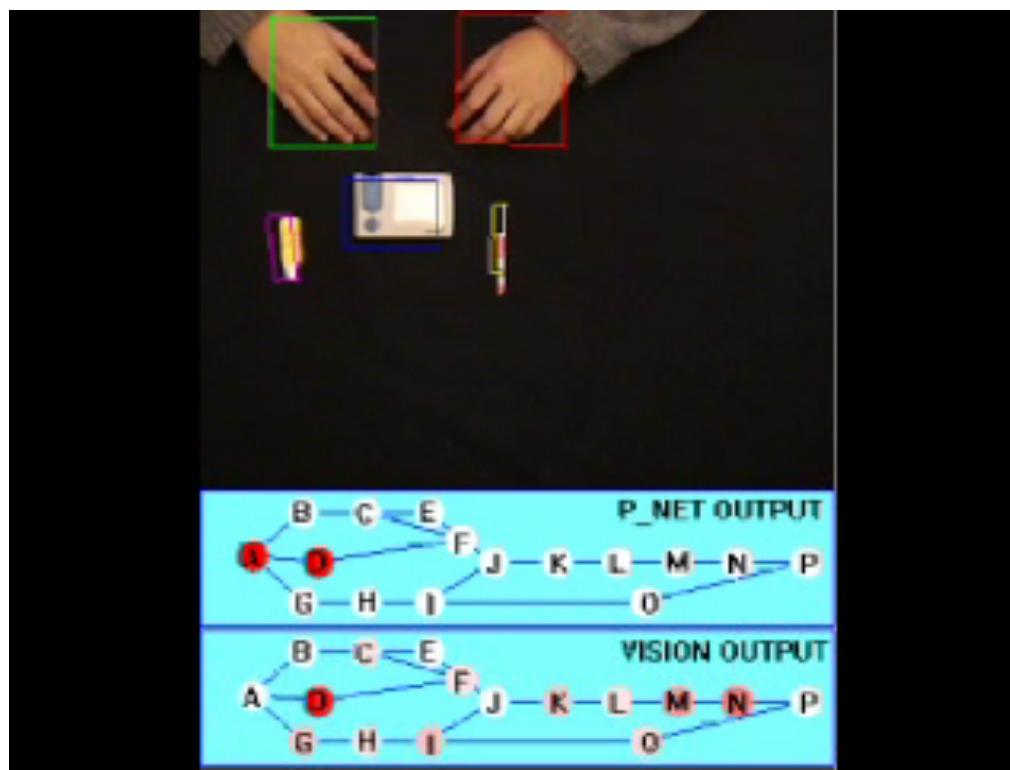
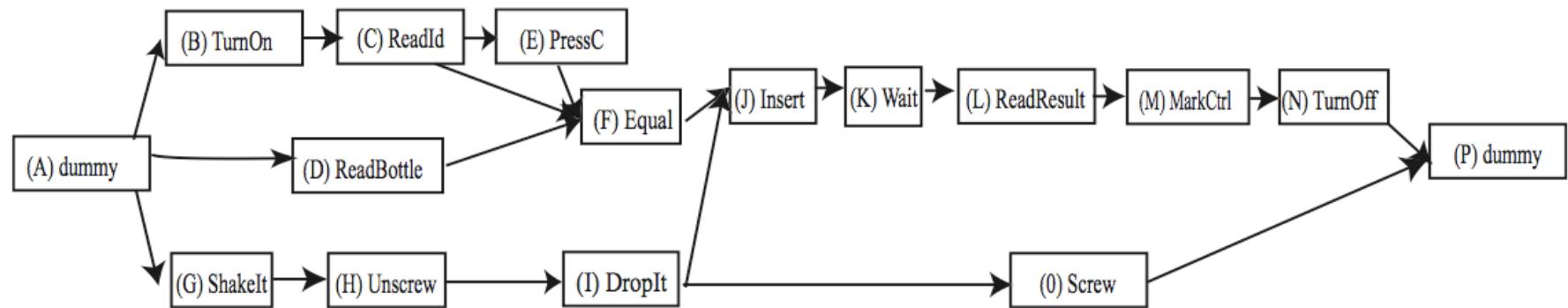
**Solution:**  
“**stand-up**” model for  
concurrent activities



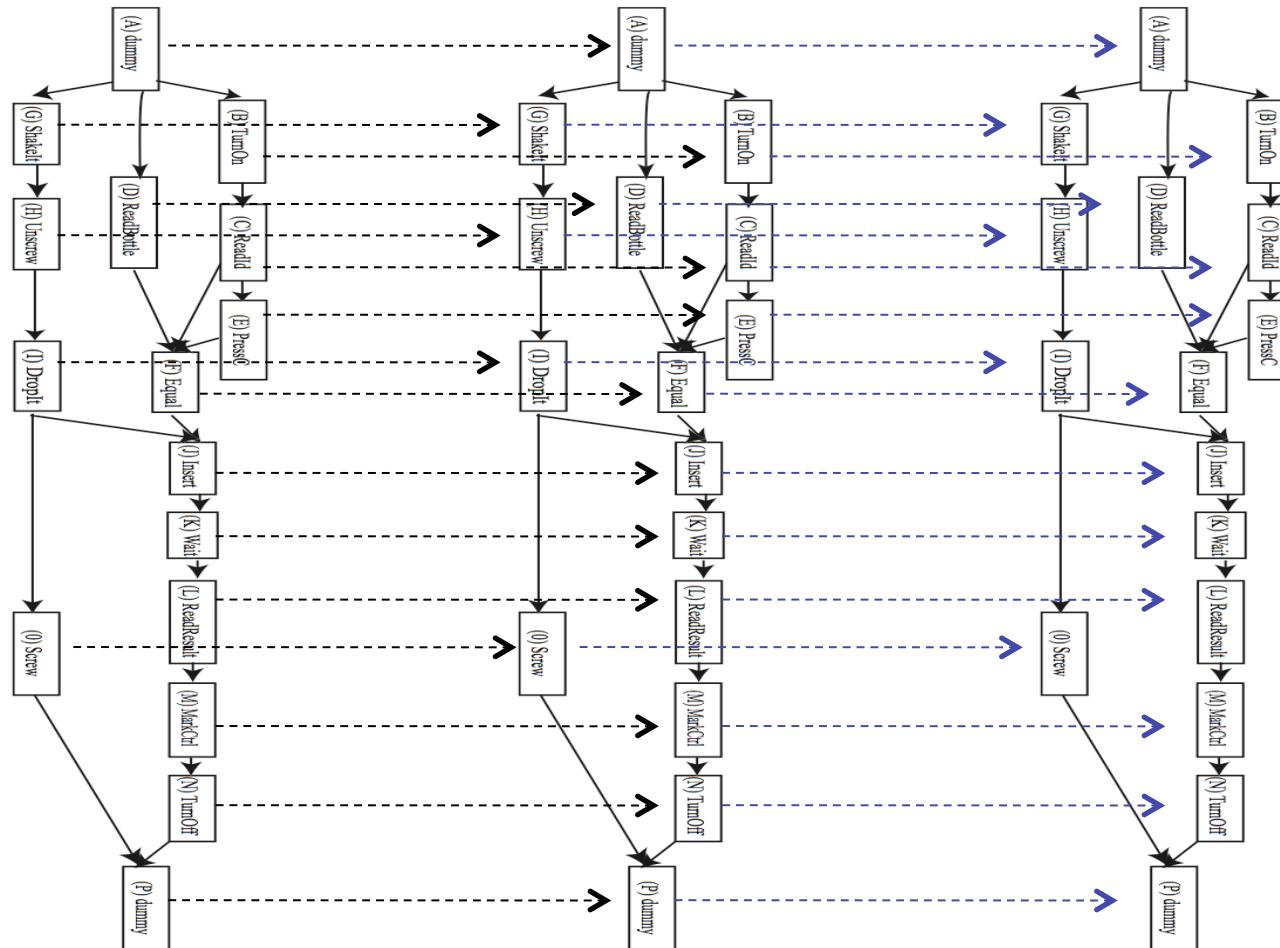
# Propagation network



Propagation Networks for Recognition of Partially Ordered Sequential Action. Shi et al 2004



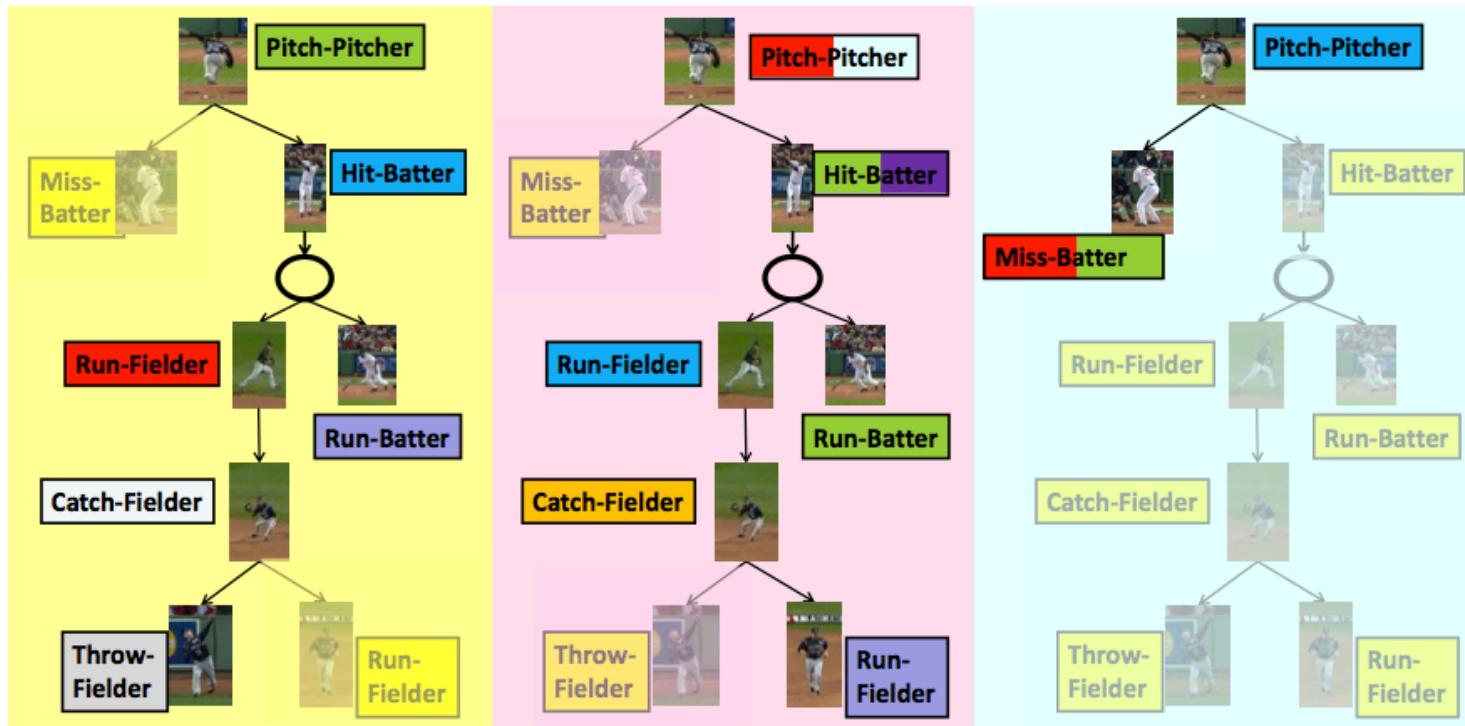
# temporal inference



Inference by standing the state transition model on its side

# Inferring structure (storylines)

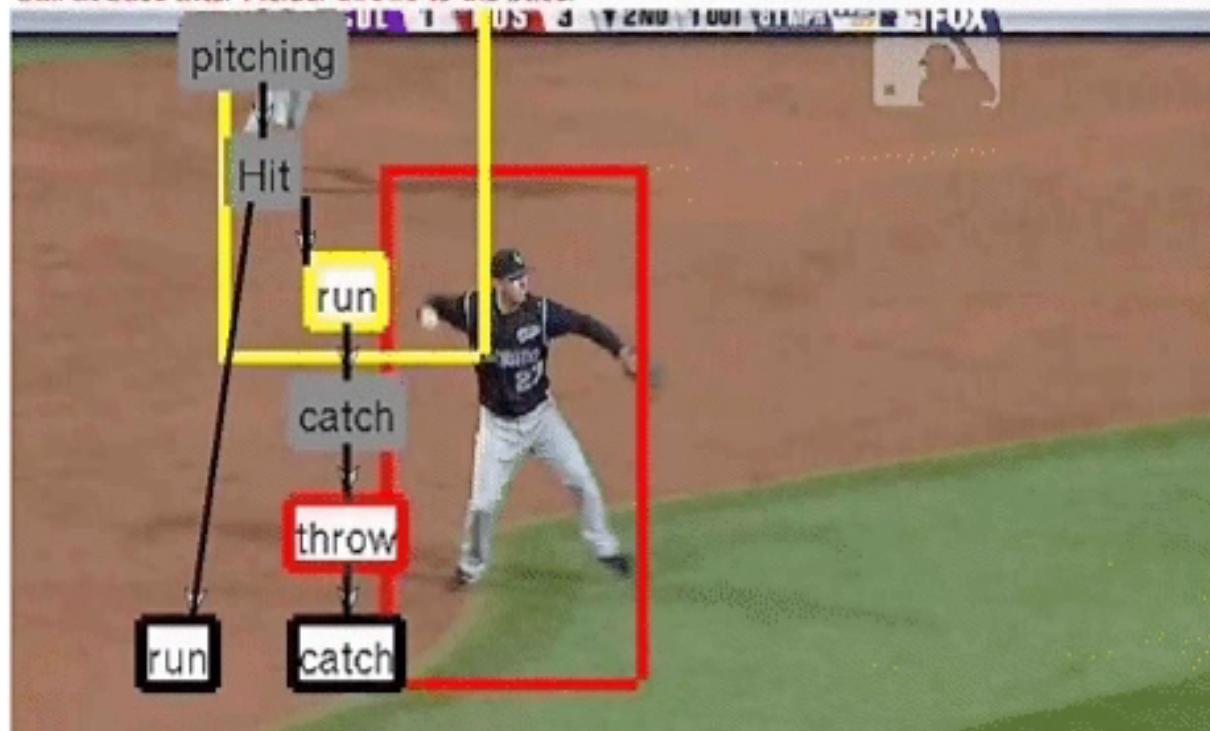
Understanding Videos, Constructing Plots –  
Learning a Visually Grounded Storyline Model from Annotated Videos  
Gupta, Srinivasan, Shi and Davis CVPR 2009



Learn AND-OR graphs from weakly labeled data

# Scripts from structure

Pitcher pitches the ball before Batter hits. Batter hits and then simultaneously Batter runs to base and Fielder runs towards the ball. Fielder runs towards the ball and then Fielder catches the ball. Fielder catches the ball and then Fielder throws to the base. Fielder at Base catches the ball at base after Fielder throws to the base.



# Take home message

Hierarchical statistical model

---

- Use when
  - Low-level action detectors are noisy
  - Structure of activity is sequential
  - Integrating dynamics
- Not for
  - Activities with deep hierarchical structure
  - Activities with complex temporal structure

# Contrasting hierarchical approaches

	<b>Actions as:</b>	<b>Activities as:</b>	<b>Model</b>	<b>Characteristic</b>
Statistic	probabilistic states	paths	DBN	Robust to uncertainty
Syntactic	discrete symbols	strings	CFG	Describes deep hierarchy
Descriptive	logical relationships	sets	CFG, MLN	Encodes complex logic

# References

(not included in ACM survey paper)

---

- W. Tsai and K.S. Fu. Attributed Grammar-A Tool for Combining Syntactic and Statistical Approaches to Pattern Recognition. SMC1980.
- M. Brand. The "Inverse Hollywood Problem": From video to scripts and storyboards via causal analysis. AAAI 1997.
- T. Wang, H. Shum, Y. Xu, N. Zheng. Unsupervised Analysis of Human Gestures. PRCM 2001.
- C.G. Nevill-Manning, I.H. Witten. On-Line and Off-Line Heuristics for Inferring Hierarchies of Repetitions in Sequences. IEEE 2000.
- K. Heller, Y.W. Teh and D. Gorur. Infinite Hierarchical Hidden Markov Model s. AISTATS 2009.
- P. Liang, S. Petrov, M. Jordan, D. Klein. The Infinite PCFG using Hierarchical Dirichlet Processes. EMNLP 2007.
- A. Gupta, N. Srinivasan, J. Shi and L. Davis. Understanding Videos, Constructing Plots - Learning a Visually Grounded Storyline Model from Annotated Videos. CVPR 2009.