

Dialogue Semantics and Pragmatics

A Tutorial at the ESSENCE Fall School 2014
Part II

David Schlangen
Universität Bielefeld, Germany
david.schlangen@uni-bielefeld.de

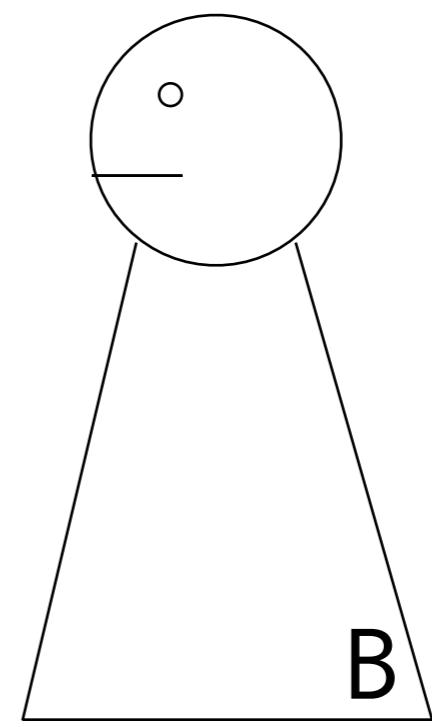
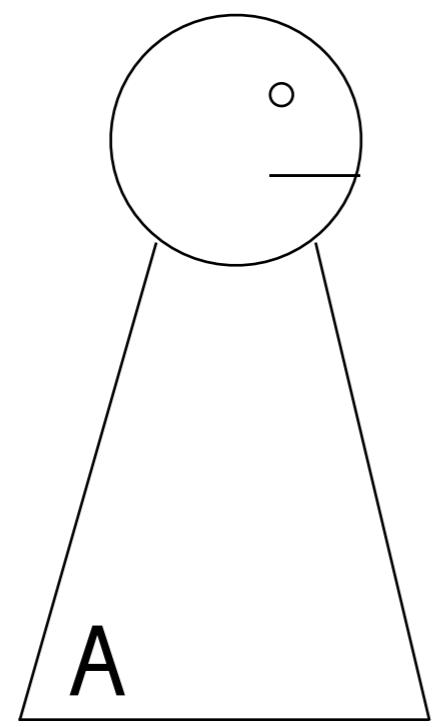
Overview

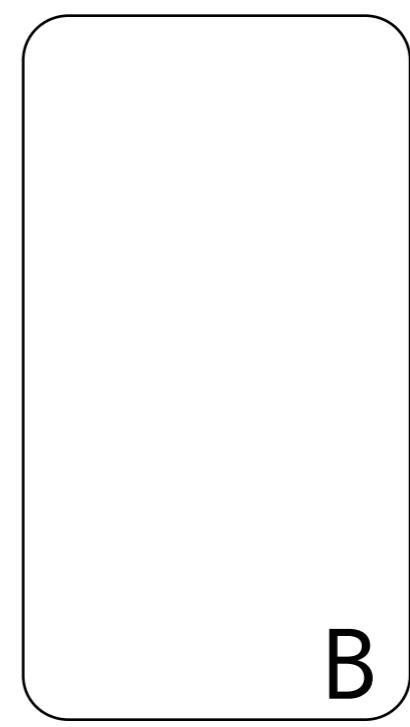
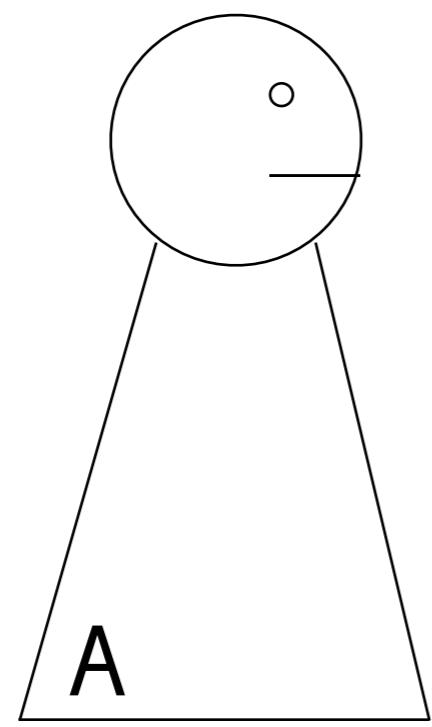
- **Part I: Foundations**

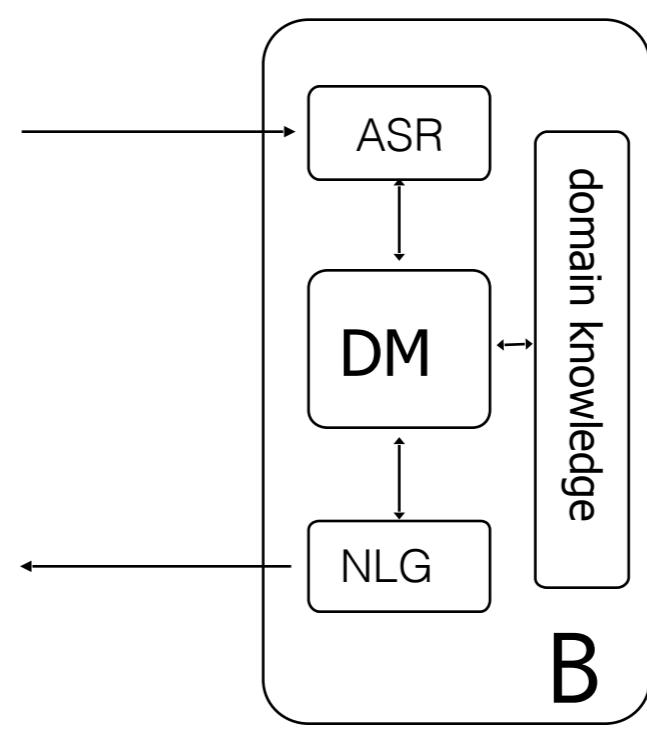
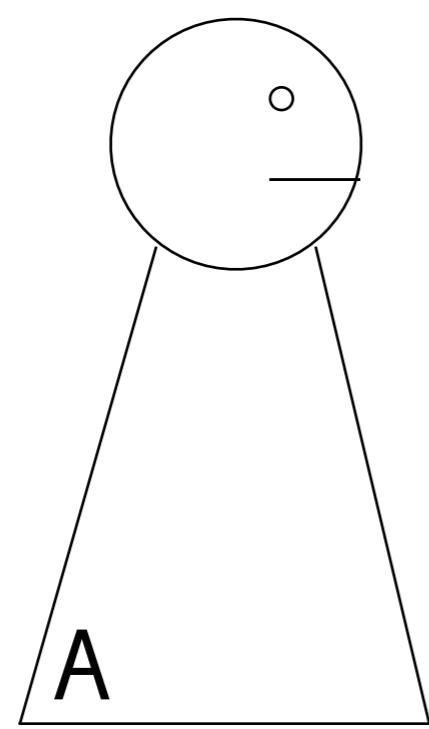
- coordination, convention
- communicative intentions
- non-conventional meaning
- grounding
- turn-taking
- disfluencies

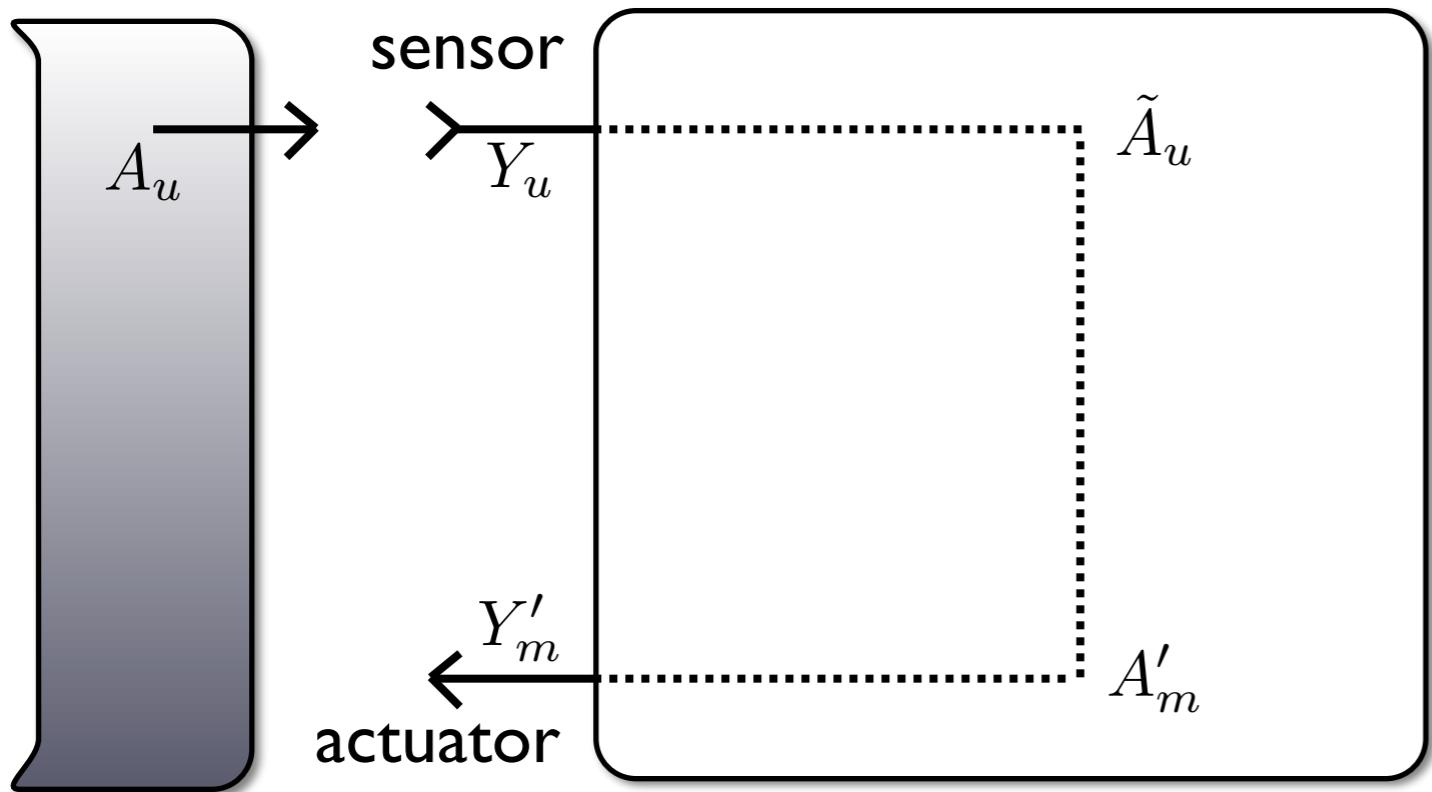
- **Part II: Computational Models**

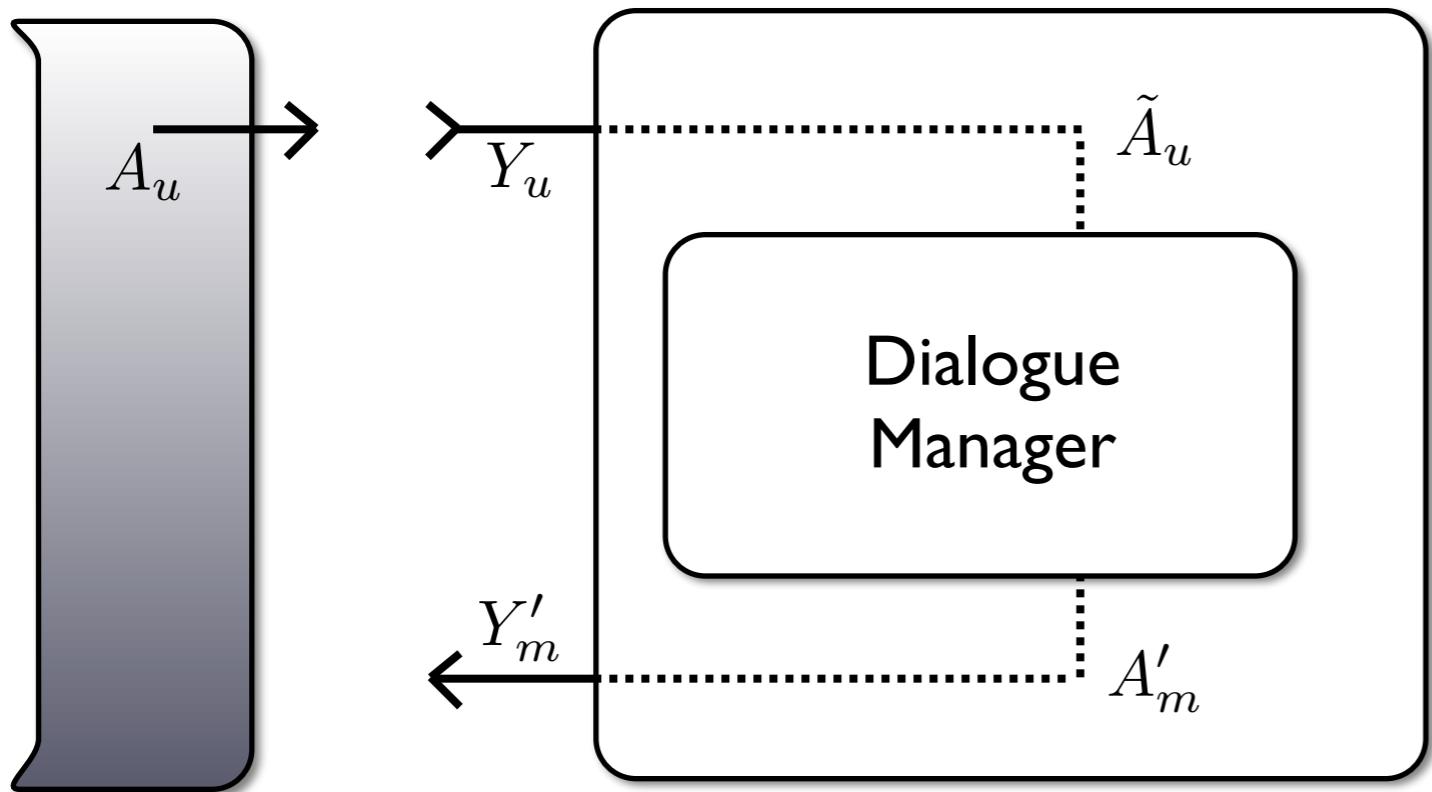
- approaches to dialogue modelling
- incremental processing, turn-taking
- an example: grounded semantics

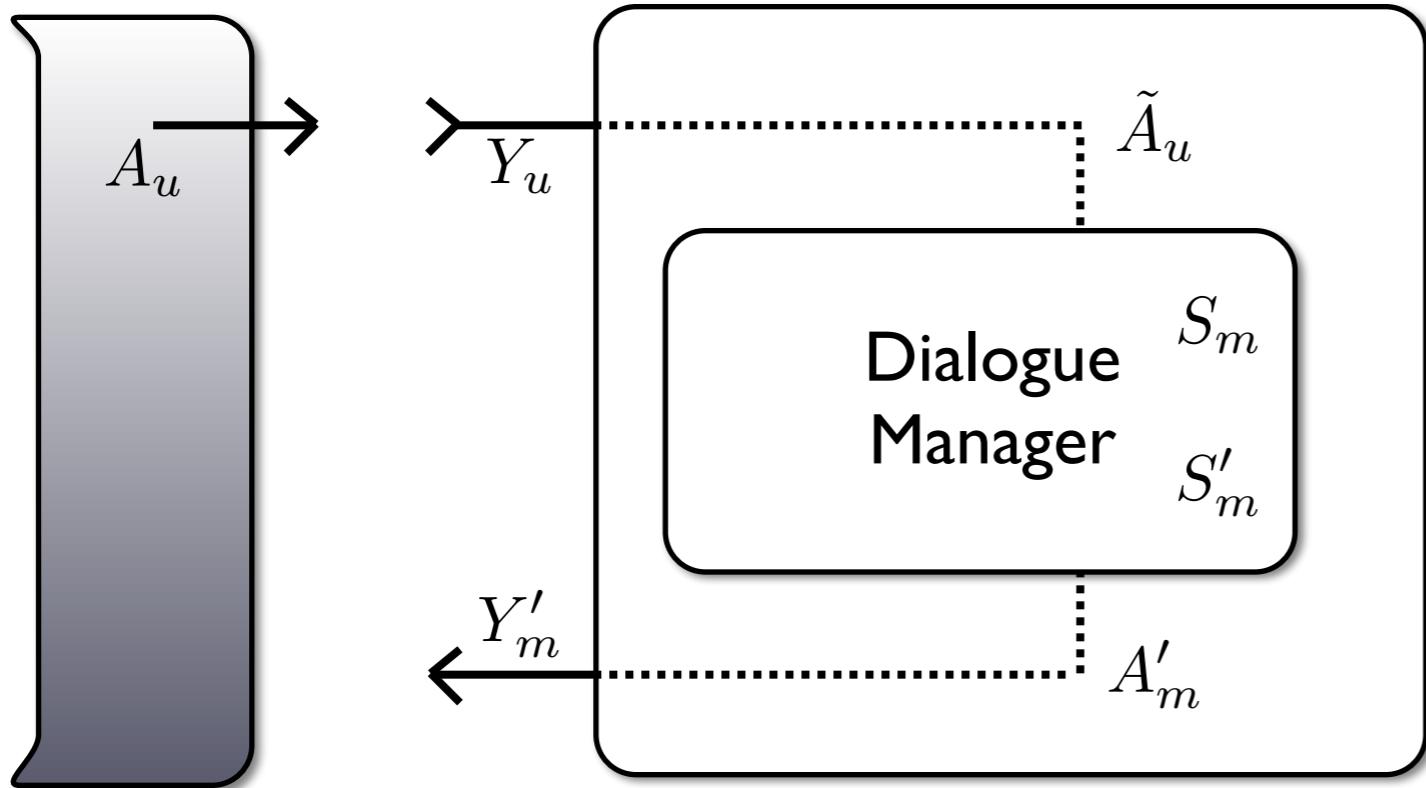








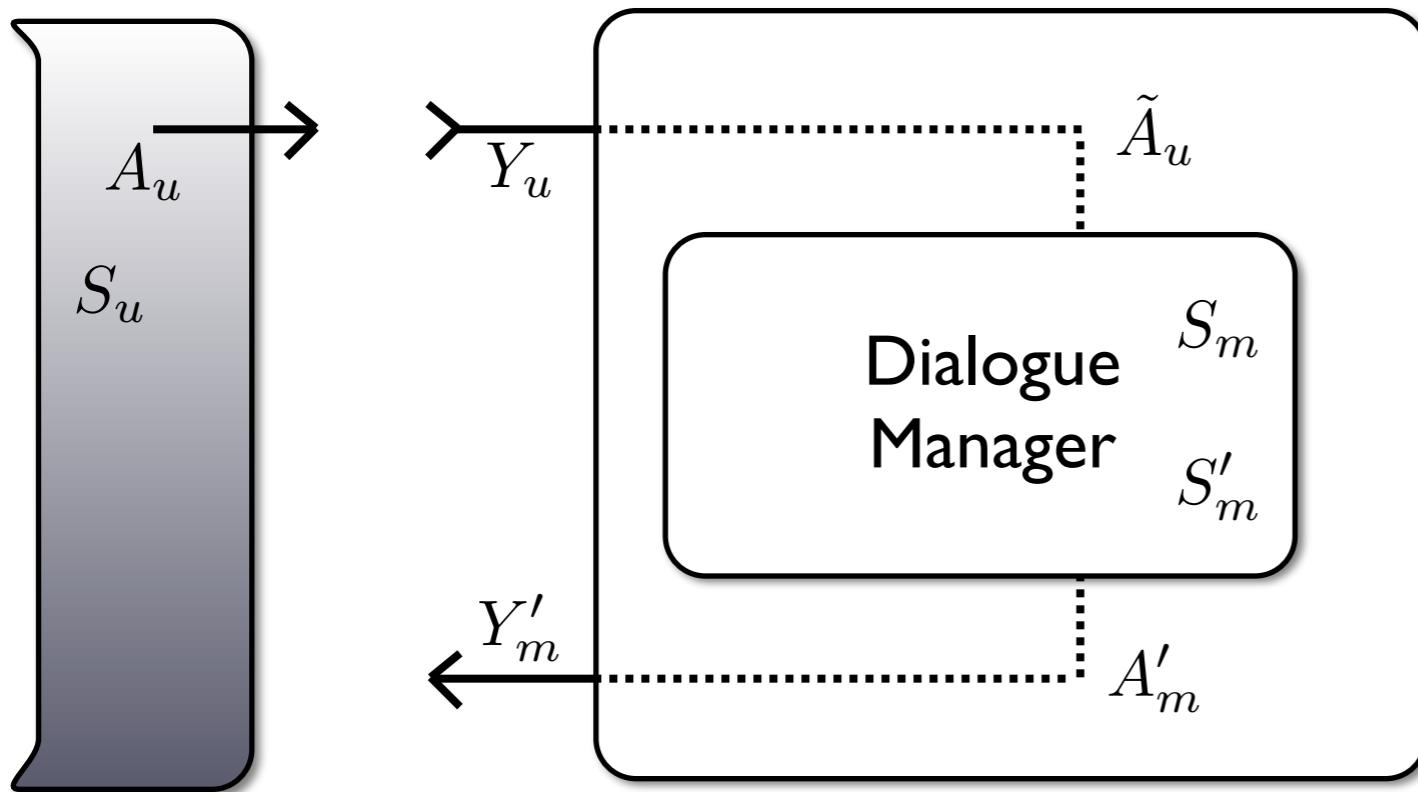




$$U(S_m, \tilde{A}_u) = S'_m$$

$$G(S_m) = A_m$$

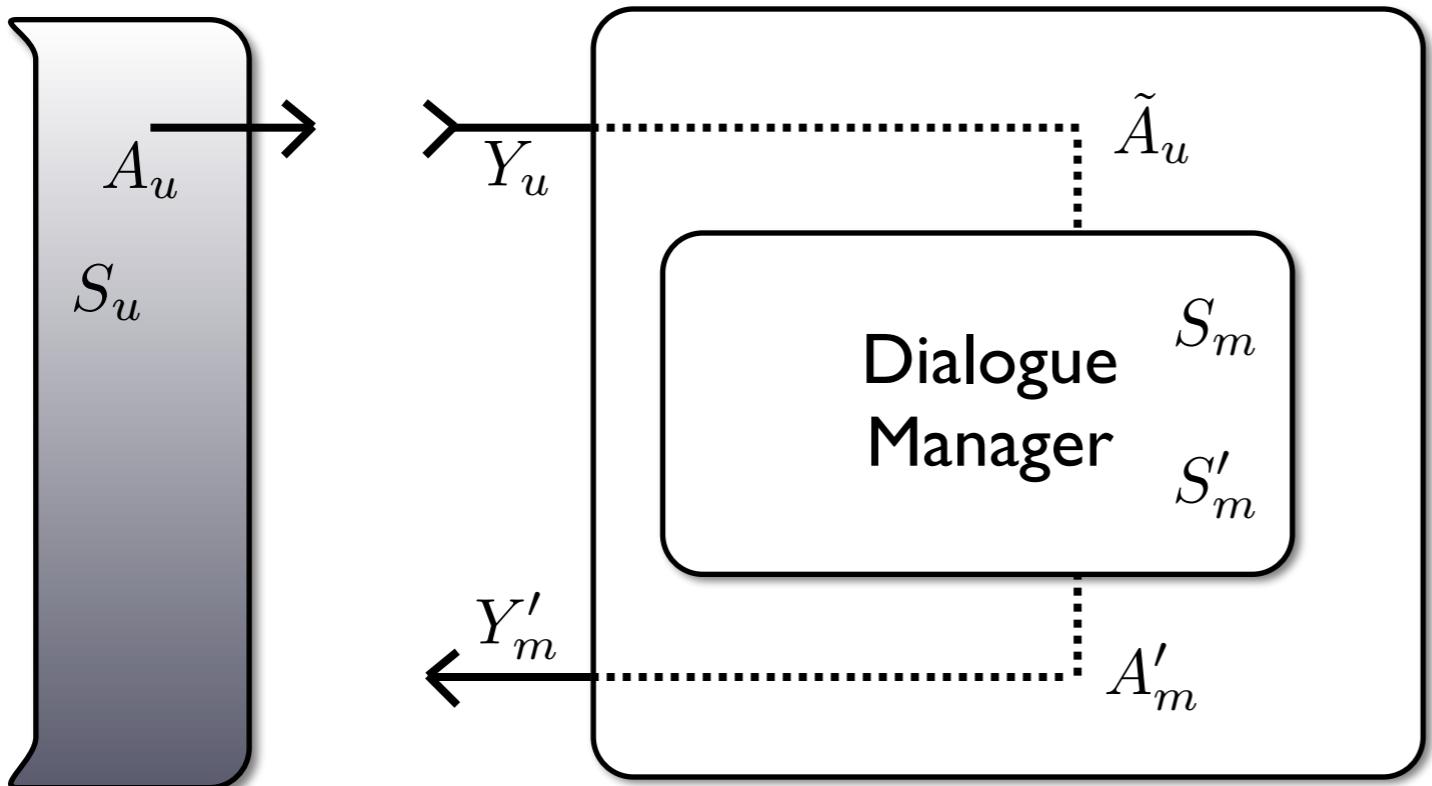
- **FSA** (eg. Oviatt et al. 1994, McTear 1998)
- **Forms** (eg., VoiceXML; Goddeau et al. 1996)
- **Stacks, Trees, ...** (Bohus & Rudnicky 2009; Lemon et al. 2002)
- -- **(PO)MDP** (Singh et al. 2000; Williams & Young 2007)
- **Conversational Scoreboard** (Larsson & Traum 2000) (Schlangen 2003)
- **BDI** (Allen 1995; Perrault & Allen 1980)



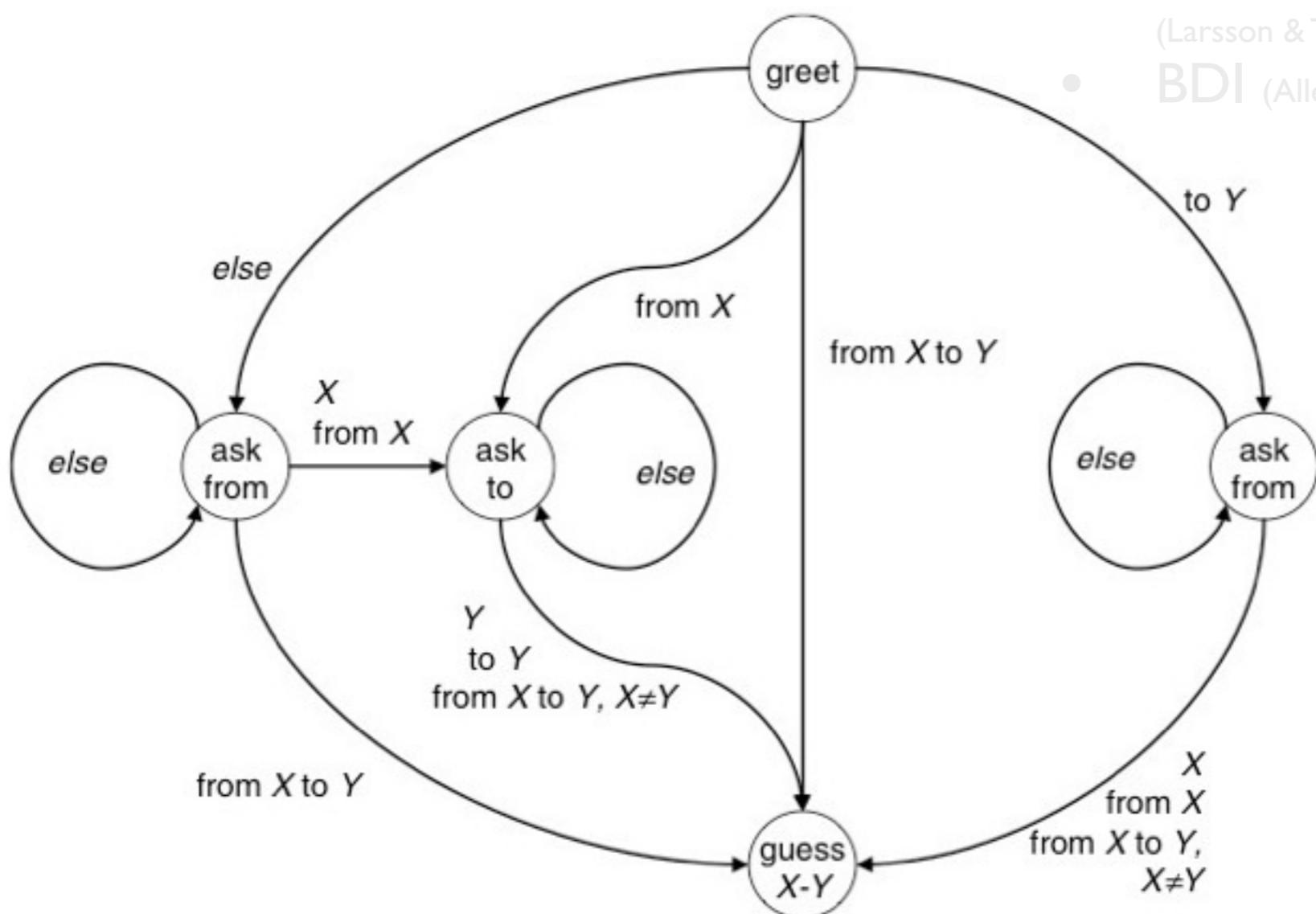
$$U(S_m, \tilde{A}_u) = S'_m$$

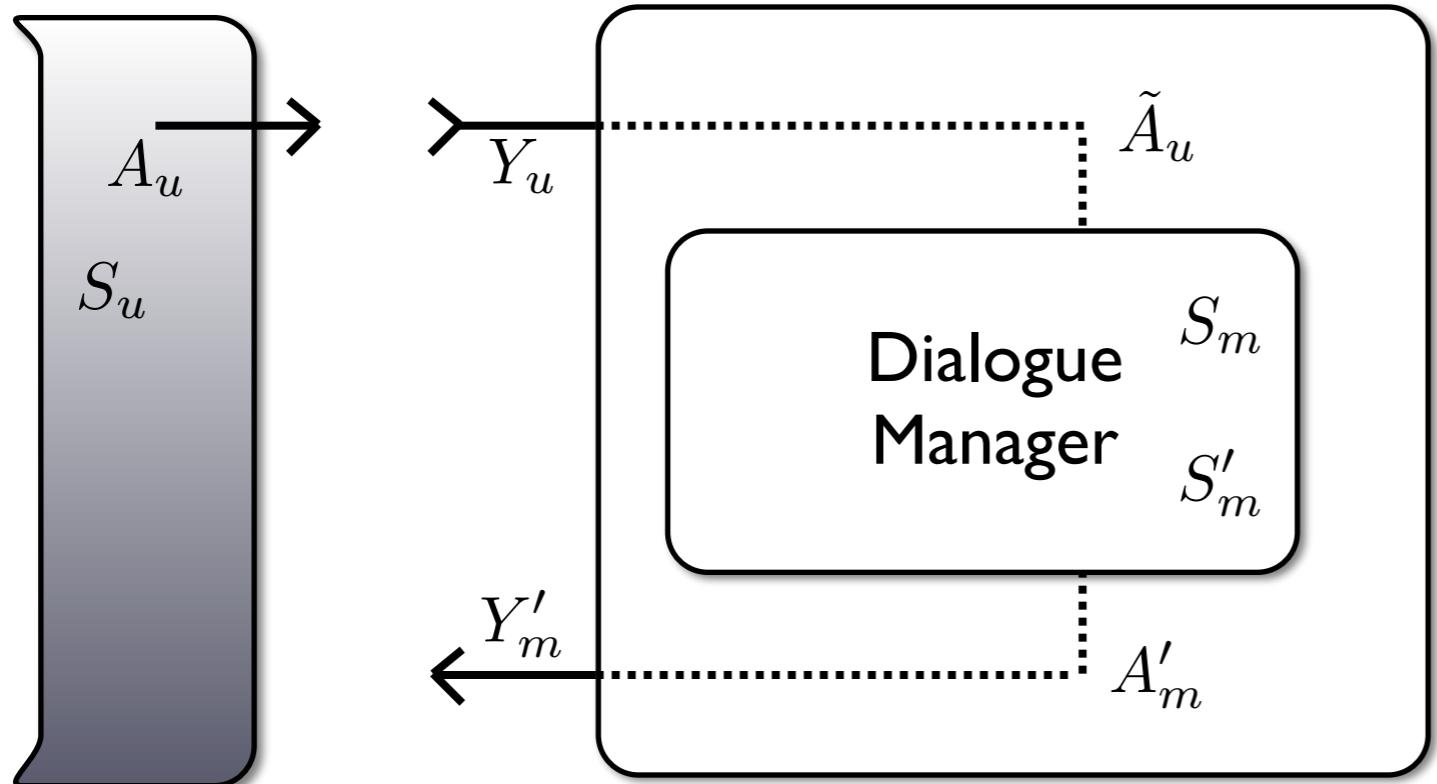
$$G(S_m) = A_m$$

- **FSA** (eg. Oviatt et al. 1994, McTear 1998)
- **Forms** (eg., VoiceXML; Goddeau et al. 1996)
- **Stacks, Trees, ...** (Bohus & Rudnicky 2009; Lemon et al. 2002)
- **-- (PO)MDP** (Singh et al. 2000; Williams & Young 2007)
- **Conversational Scoreboard** (Larsson & Traum 2000)(Schlangen 2003)
- **BDI** (Allen 1995; Perrault & Allen 1980)



- **FSA** (eg. Oviatt et al. 1994, McTear 1998)
- **Forms** (eg., VoiceXML; Goddeau et al. 1996)
- **Stacks, Trees, ...** (Bohus & Rudnicky 2009; Lemon et al. 2002)
- **-- (PO)MDP** (Singh et al. 2000; Williams & Young 2007)
- **Conversational Scoreboard** (Larsson & Traum 2000)(Schlangen 2003)
- **BDI** (Allen 1995; Perrault & Allen 1980)





State: Beliefs, Desires, Intentions

Update: inference rules

Decisions: planning

“Can you give me a list of flights to Berlin?”

```

S.REQUEST(S,H,InformIf(H,S,CanDo(H,Give(H,S,List))))
B(H,W(S,InformIf(H,S,CanDo(H,Give(H,S,List)))))
B(H,W(S,KnowIf(H,S,CanDo(H,Give(H,S,List)))))
B(H,W(S,CanDo(H,Give(H,S,List))))
B(H,W(S,Give(H,S,List)))
REQUEST(H,S,Give(H,S,List))
W(H,Give(H,S,List))

```

W(S,Give(S,H,List))

Action: Inform(S,H,P)

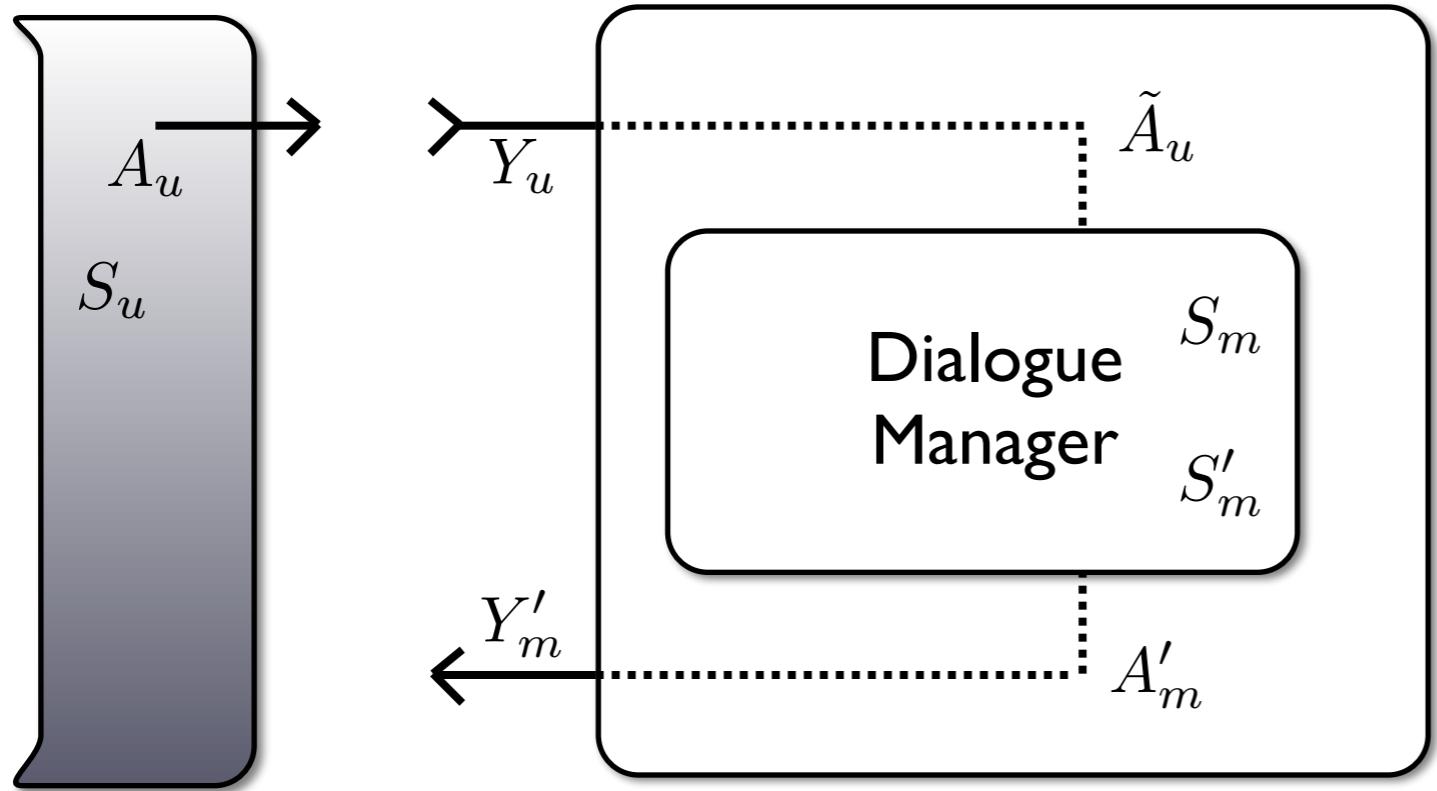
Effect: Know(H,P)

$$U(S_m, \tilde{A}_u) = S'_m$$

$$G(S_m) = A_m$$

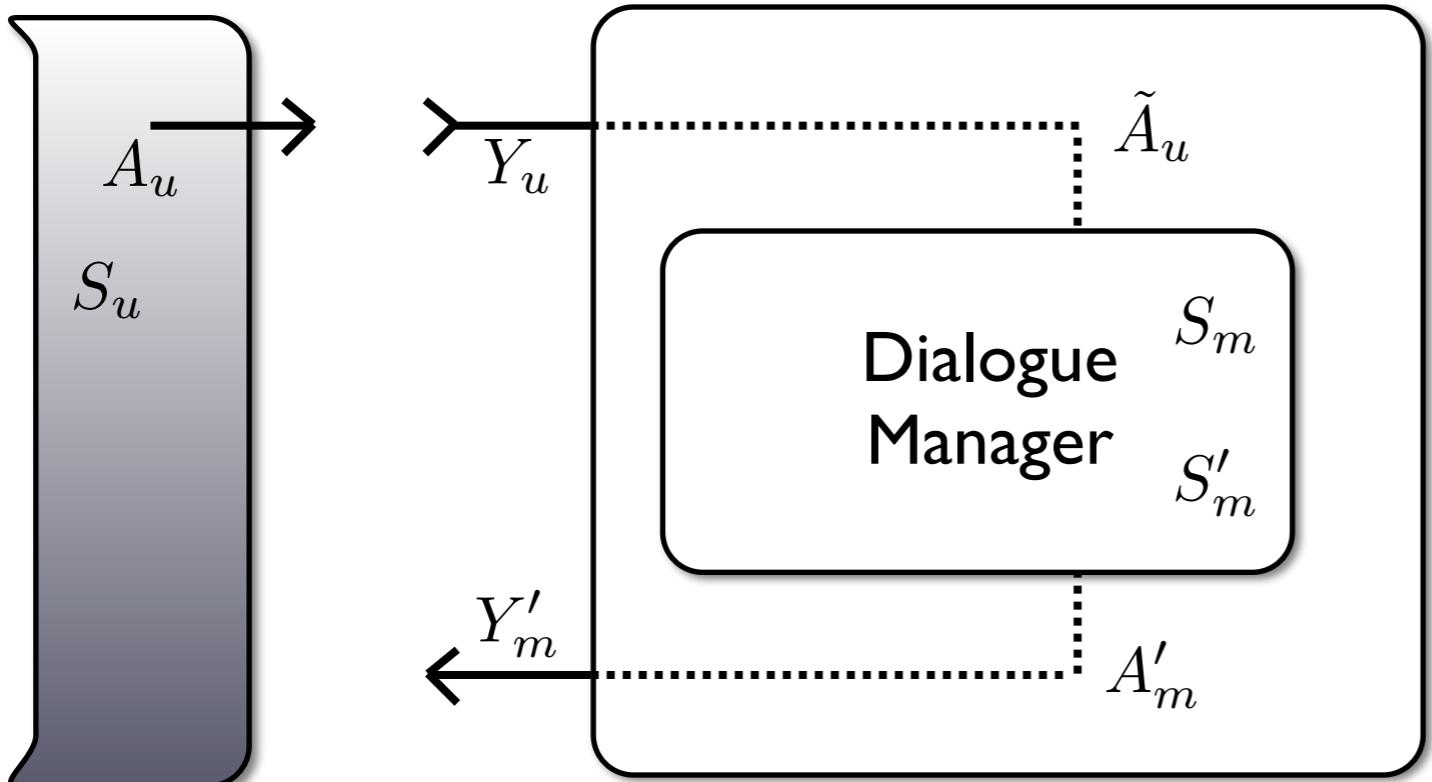
- FSA (eg. Oviatt et al. 1994, McTear 1998)
- Forms (eg., VoiceXML; Goddeau et al. 1996)
- Stacks, Trees, ... (Bohus & Rudnicky 2009; Lemon et al. 2002)
- -- (PO)MDP (Singh et al. 2000; Williams & Young 2007)
- Conversational Scoreboard (Larsson & Traum 2000) (Schlangen 2003)
- BDI (Allen 1995; Perrault & Allen 1980)

“There are three flights, one at ...”



- $$U(S_m, \tilde{A}_u) = S'_m$$
- $$G(S_m) = A_m$$
- **FSA** (eg. Oviatt et al. 1994, McTear 1998)
 - **Forms** (eg., VoiceXML; Goddeau et al. 1996)
 - **Stacks, Trees, ...** (Bohus & Rudnicky 2009; Lemon et al. 2002)
 - -- **(PO)MDP** (Singh et al. 2000; Williams & Young 2007)
 - **Conversational Scoreboard**
(Larsson & Traum 2000)
 - **BDI** (Allen 1995; Perrault & Allen 1980)

PRIVATE :	$\left[\begin{array}{l} \text{AGENDA} : \text{Stack(Action)} \\ \text{PLAN} : \text{Stack(Action)} \\ \text{BEL} : \text{Set(Prop)} \end{array} \right]$
SHARED :	$\left[\begin{array}{l} \text{COM} : \text{Set(Prop)} \\ \text{QUD} : \text{Stack(Question)} \\ \text{LU} : \left[\begin{array}{l} \text{SPEAKER} : \text{Participant} \\ \text{MOVE} : \text{Move} \end{array} \right] \end{array} \right]$



"price information,
please"

```

getLatestMove
{
    set(/SHARED/LU/MOVES, set([ask(?A.price(A))]))
    set(/SHARED/LU/SPEAKER, usr)
}
integrateUsrAsk
{
    push(/SHARED/QUD, ?A.price(A))
    push(/PRIVATE/AGENDA, respond(?A.price(A)))
}
findPlan
{
    pop(/PRIVATE/AGENDA)
    set(/PRIVATE/PLAN, stack([raise(?C.how(C)), findout(?D.dest_city(D)), ... ]))
}

```

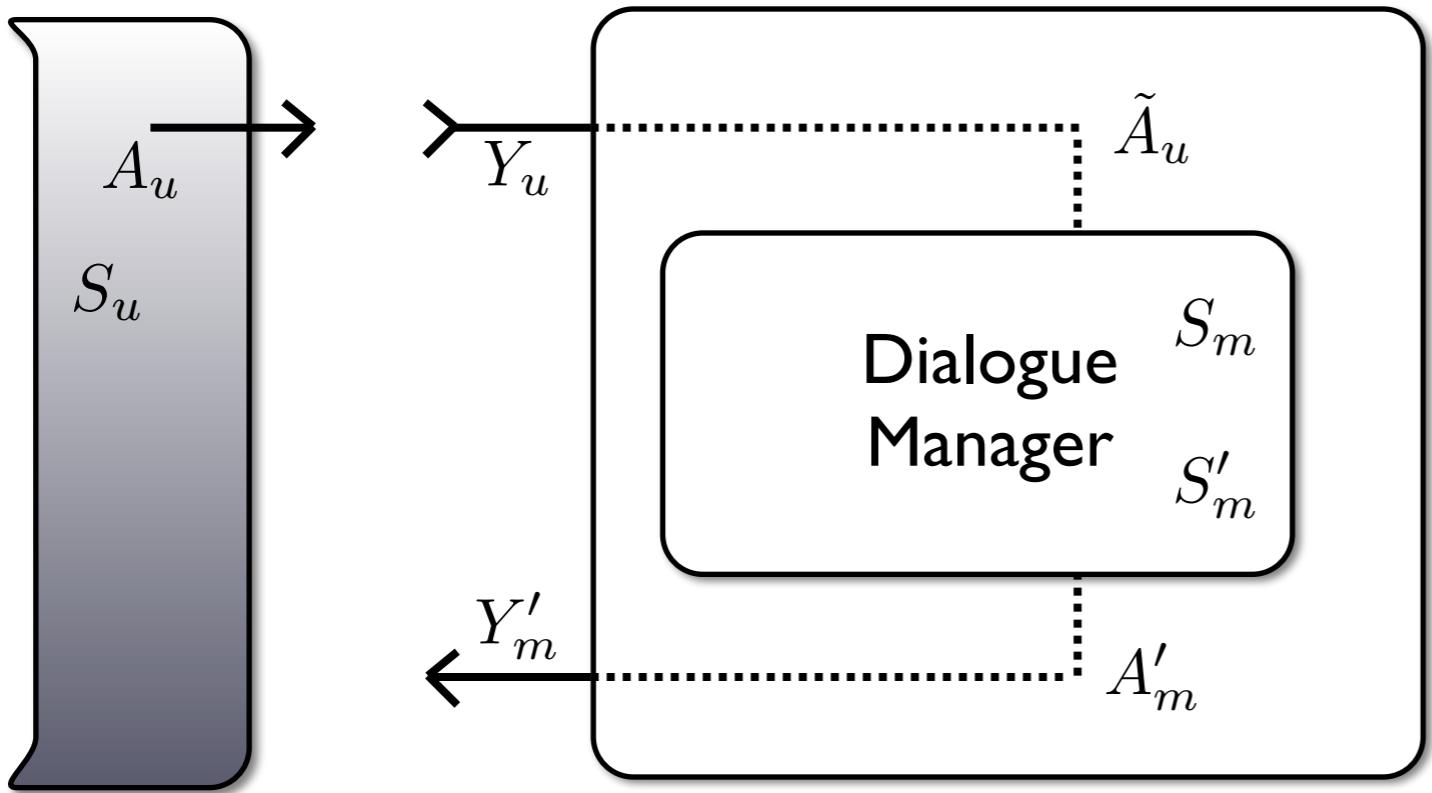
PRIVATE	AGENDA = <	<	selectFromPlan push(/PRIVATE/AGENDA, raise(?A.how(A))) selectAsk { add(NEXT_MOVES, ask(?A.how(A))) if_do(fst(\$/PRIVATE/PLAN, raise(?A.how(A)))
	raise(?A.how(A))	>	
	findout(?B.dest_city(B))		
	findout(?C.dept_city(C))		
PLAN	findout(?D.month(D))		
	findout(?E.dept_day(E))		
	findout(?F.class(F))		
	consultDB(?G.price(G))		
BEL	= { }		
	COM = { }		
	QUD = < ?H.price(H) >		
	LU = [SPEAKER = usr MOVES = { ask(?H.price(H)) }]		

$$U(S_m, \tilde{A}_u) = S'_m$$

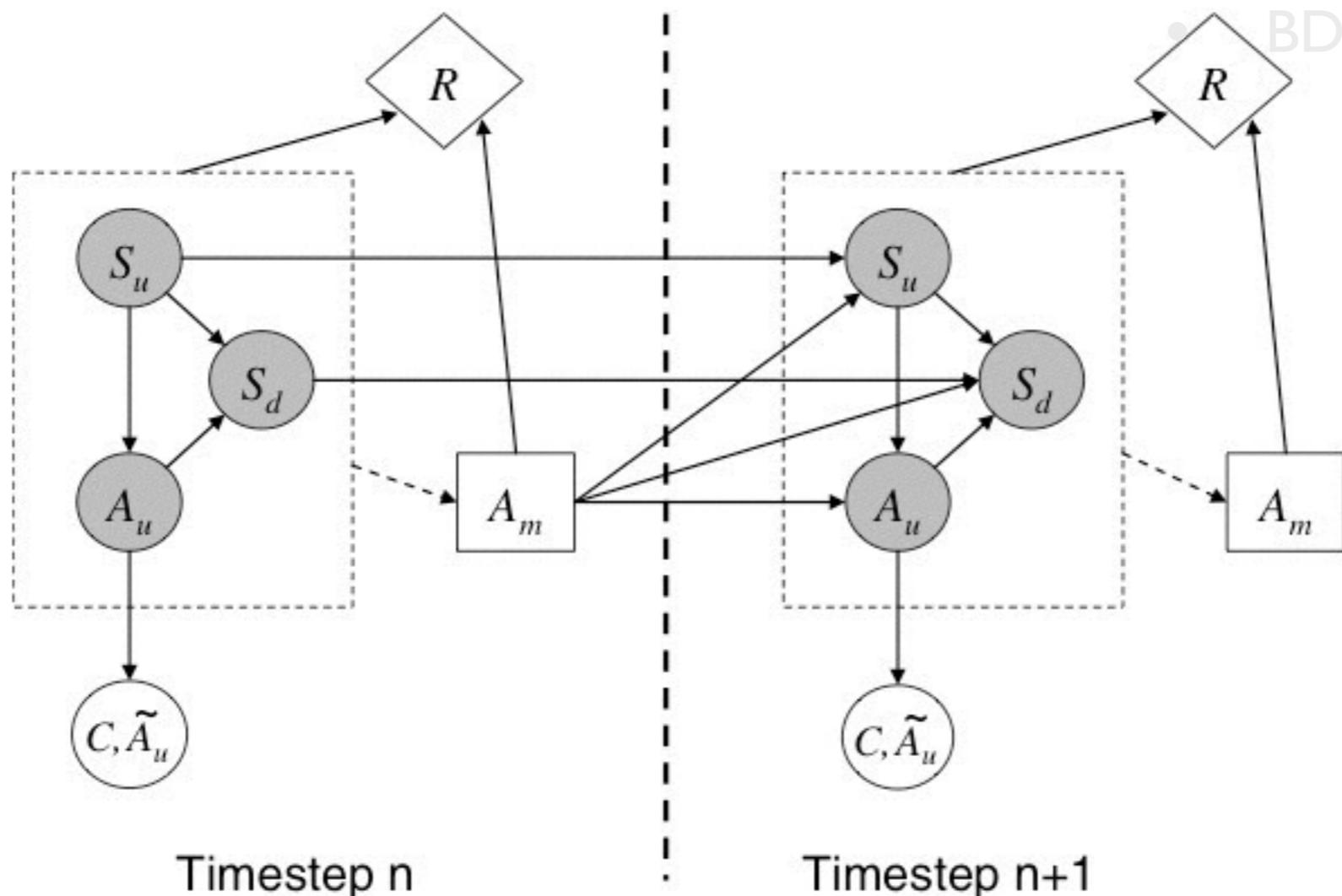
$$G(S_m) = A_m$$

- FSA (eg. Oviatt et al. 1994, McTear 1998)
- Forms (eg., VoiceXML; Goddeau et al. 1996)
- Stacks, Trees, ... (Bohus & Rudnicky 2009; Lemon et al. 2002)
- **Conversational Scoreboard**
(Larsson & Traum 2000)
- RDI (Allen 1985, Pratt & Allen 1980)

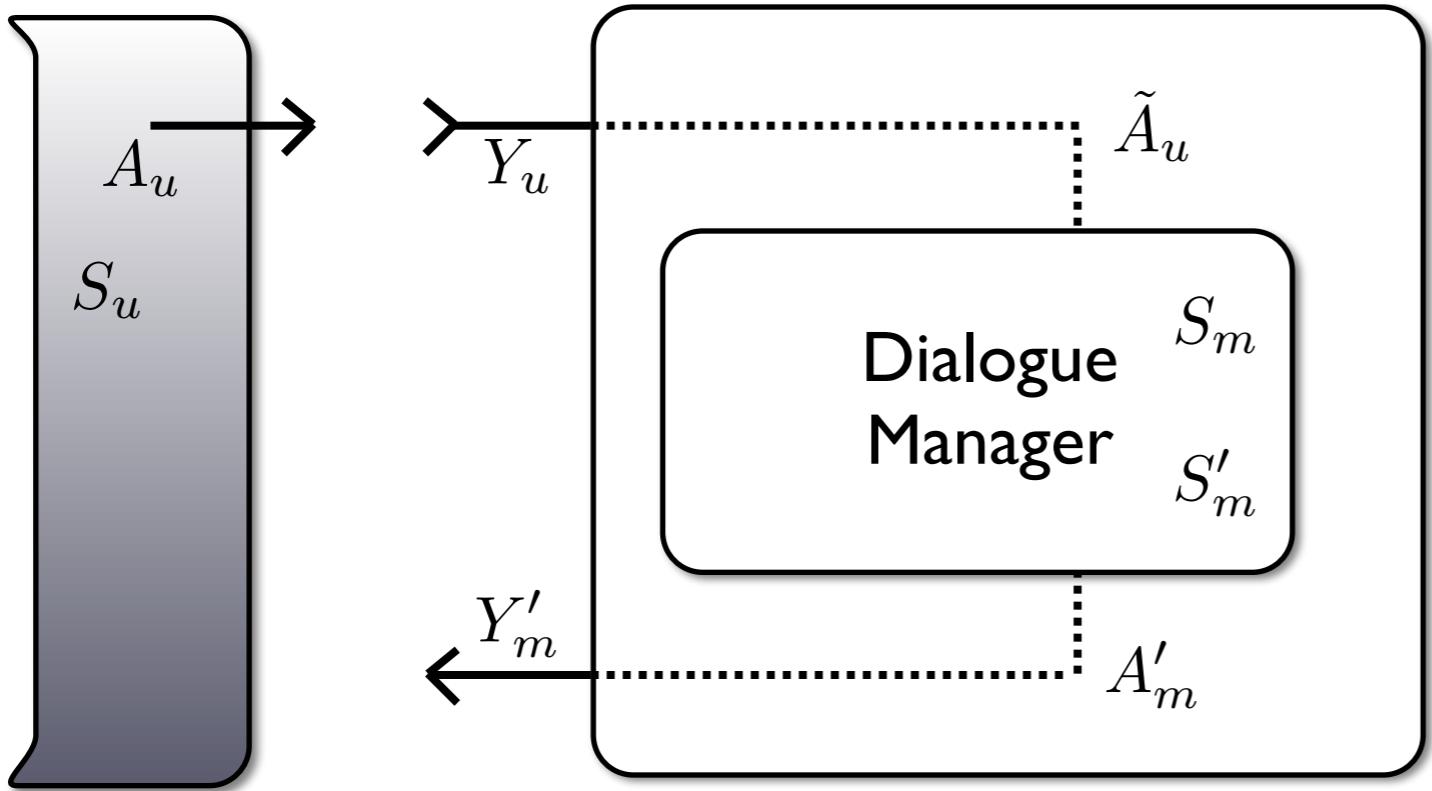
"how do you want to travel?"



- FSA (eg. Oviatt et al. 1994, McTear 1998)
- Forms (eg., VoiceXML; Goddeau et al. 1996)
- Stacks, Trees, ... (Bohus & Rudnicky 2009; Lemon et al. 2002)
- -- (PO)MDP (Singh et al. 2000; Williams & Young 2007)
- Conversational Scoreboard (Larsson & Traum 2000)(Schlangen 2003)
- BDI (Allen 1995; Perrault & Allen 1980)

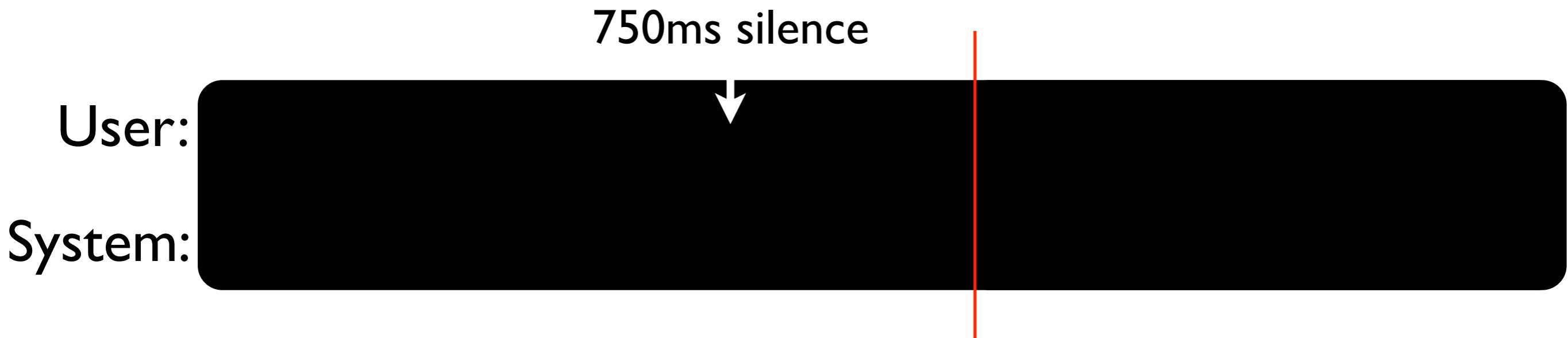


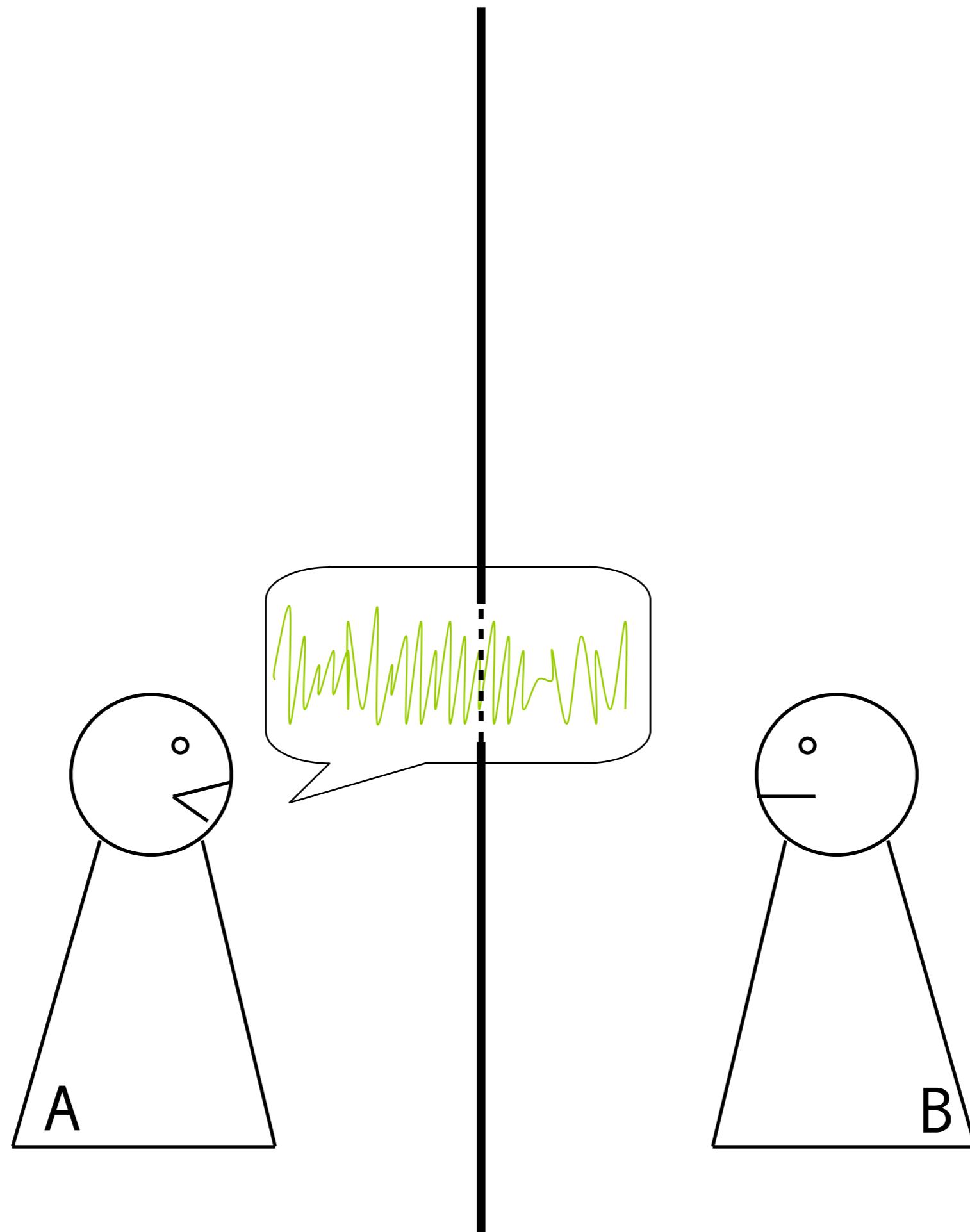
aus (Williams & Young 2007)



- $$U(S_m, \tilde{A}_u) = S'_m$$
- $$G(S_m) = A_m$$
- **FSA** (eg. Oviatt et al. 1994, McTear 1998)
 - **Forms** (eg., VoiceXML; Goddeau et al. 1996)
 - **Stacks, Trees, ...** (Bohus & Rudnicky 2009; Lemon et al. 2002)
 - **(PO)MDP** (Singh et al. 2000; Williams & Young 2007)
 - **Conversational Scoreboard**
(Larsson & Traum 2000) (Schlangen 2003)
 - **BDI** (Allen 1995; Perrault & Allen 1980)

non- Incremental Dialogue Processing



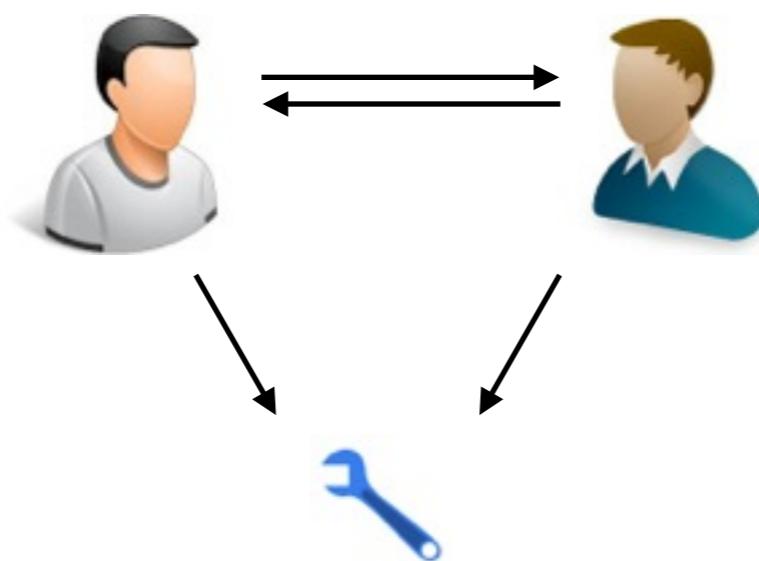


situated dialogue

- participants share a common timeline:



- participants are co-located:



Overview

- **Part I: Foundations**

- coordination, convention
- communicative intentions
- non-conventional meaning
- grounding
- turn-taking
- disfluencies

- **Part II: Computational Models**

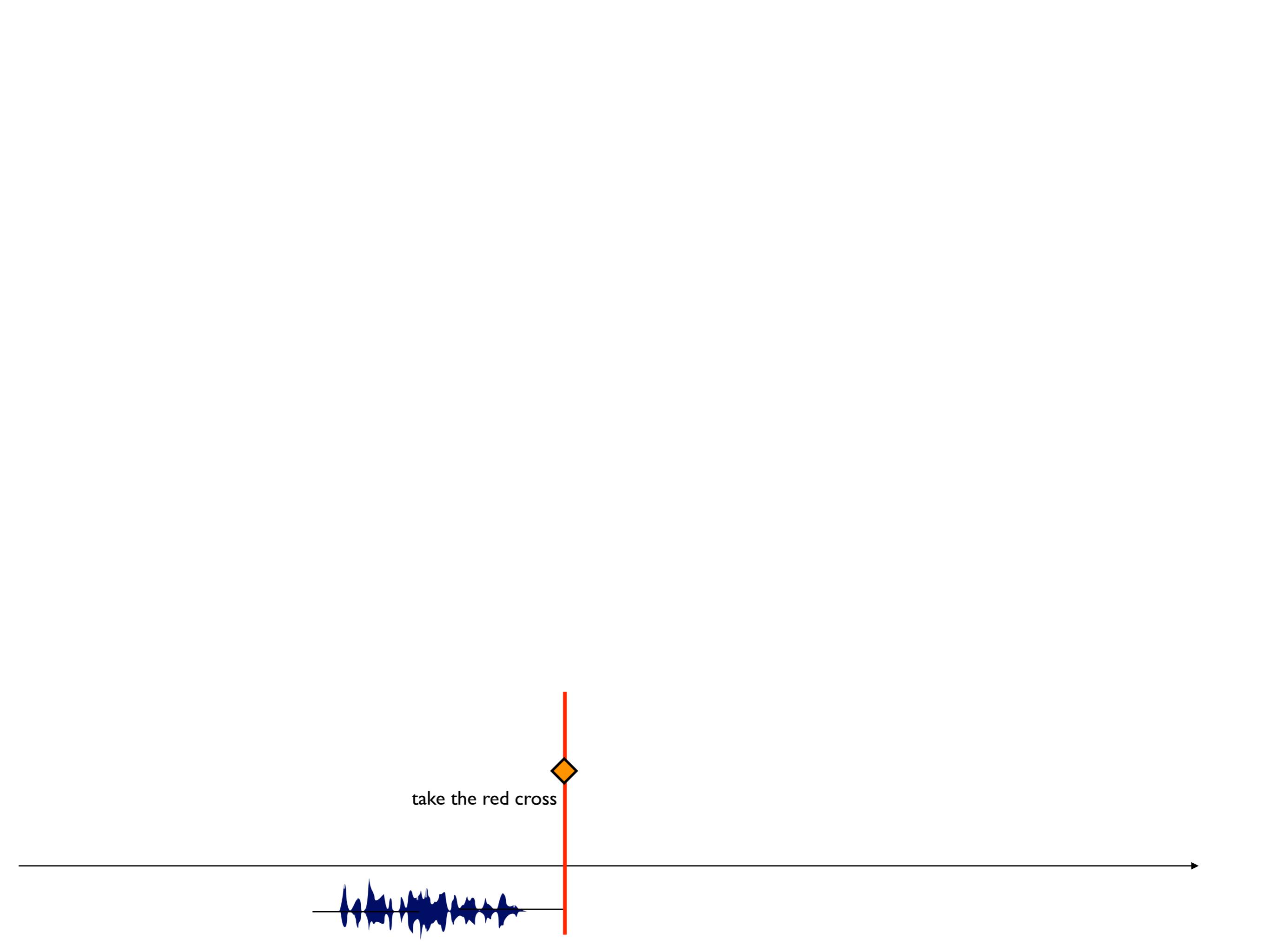
- approaches to dialogue modelling
- incremental processing, turn-taking
- an example: grounded semantics

the IU model

– Motivation –

- fundamental skills of an agent:
 - to form hypotheses about the world
 - to plan and perform actions on the world
- IU model as a temporally fine-grained model of the information state of an agent, and of how it is updated





take the red cross



the IU model

– Assumptions –

- Information state is updated with *minimal units* of information, as soon as they can be hypothesised

the IU model

– Assumptions –

- Information state is updated with *minimal units* of information, as soon as they can be hypothesised



the IU model

– Assumptions –

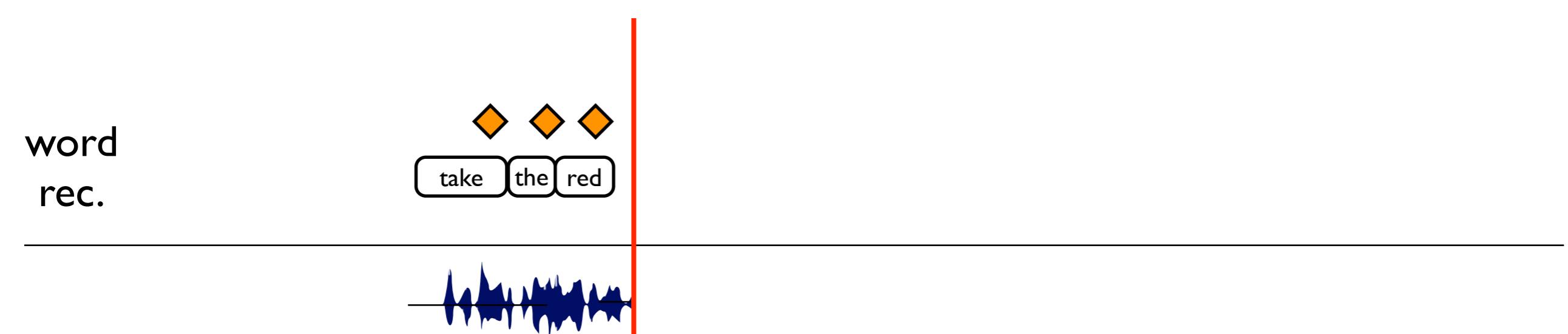
- Information state is updated with *minimal units* of information, as soon as they can be hypothesised



the IU model

– Assumptions –

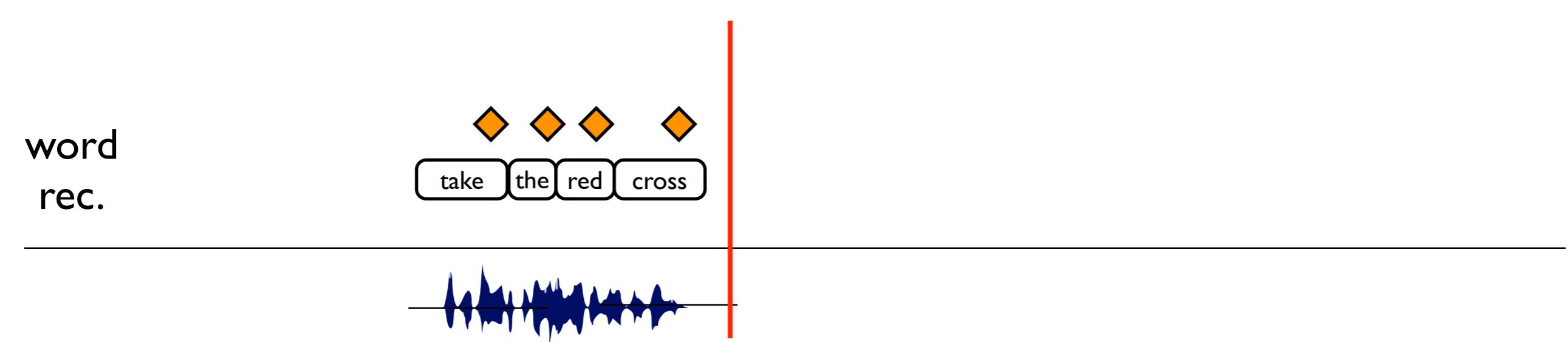
- Information state is updated with *minimal units* of information, as soon as they can be hypothesised



the IU model

– Assumptions –

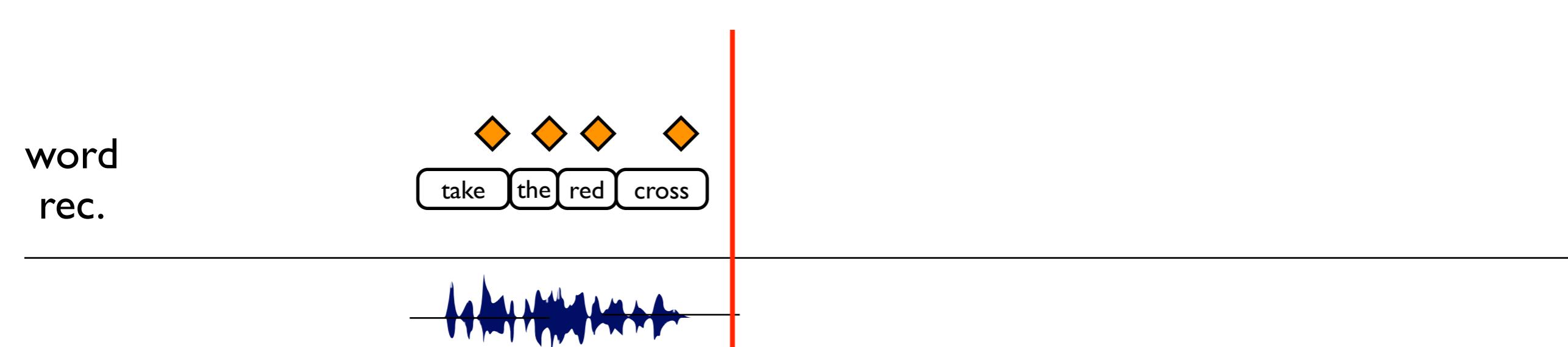
- Information state is updated with *minimal units* of information, as soon as they can be hypothesised



the IU model

– Assumptions –

- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.



the IU model

– Assumptions –

- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.



the IU model

– Assumptions –

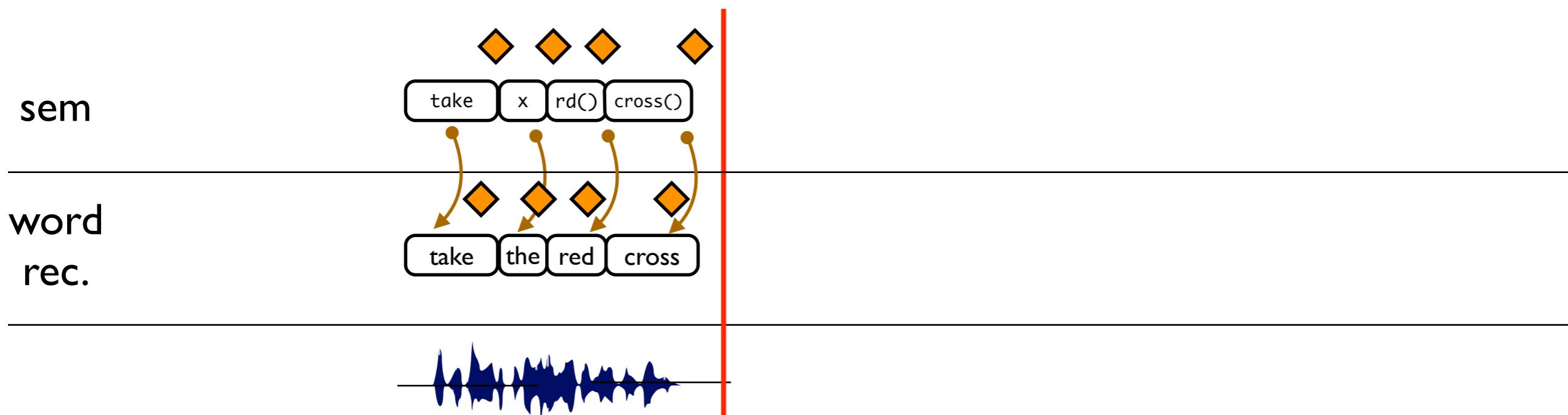
- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.



the IU model

– Assumptions –

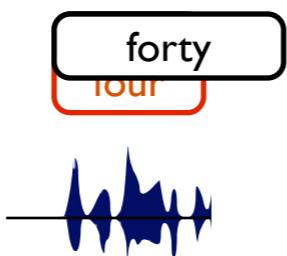
- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.



the IU model

– Assumptions –

- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.
- IS may have to be revised, in light of newer information



the IU model

– Assumptions –

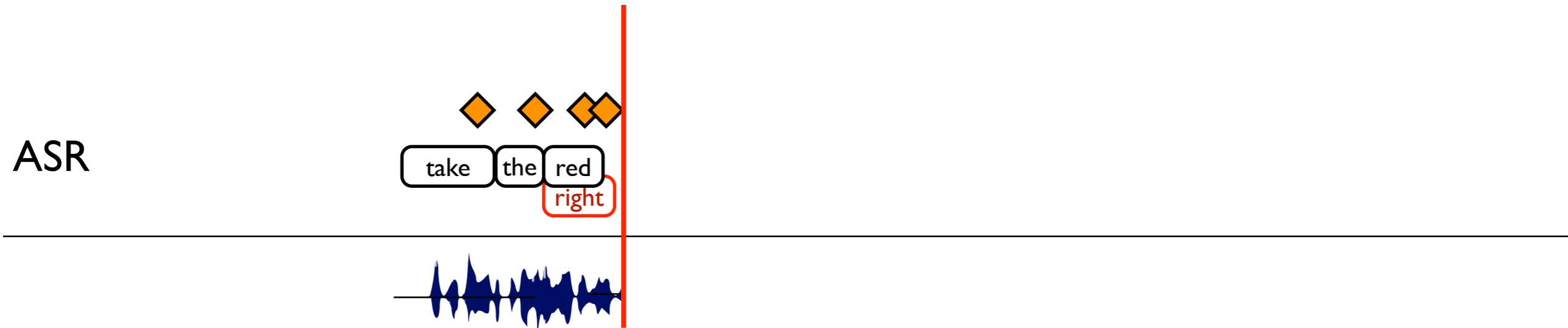
- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.
- IS may have to be revised, in light of newer information



the IU model

– Assumptions –

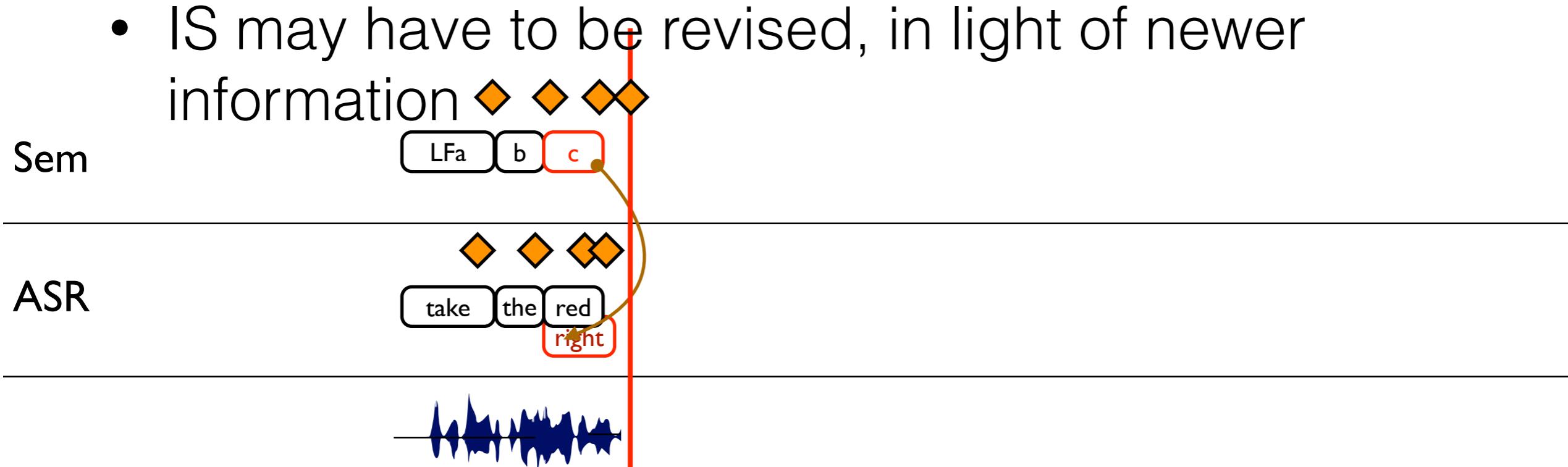
- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.
- IS may have to be revised, in light of newer information



the IU model

– Assumptions –

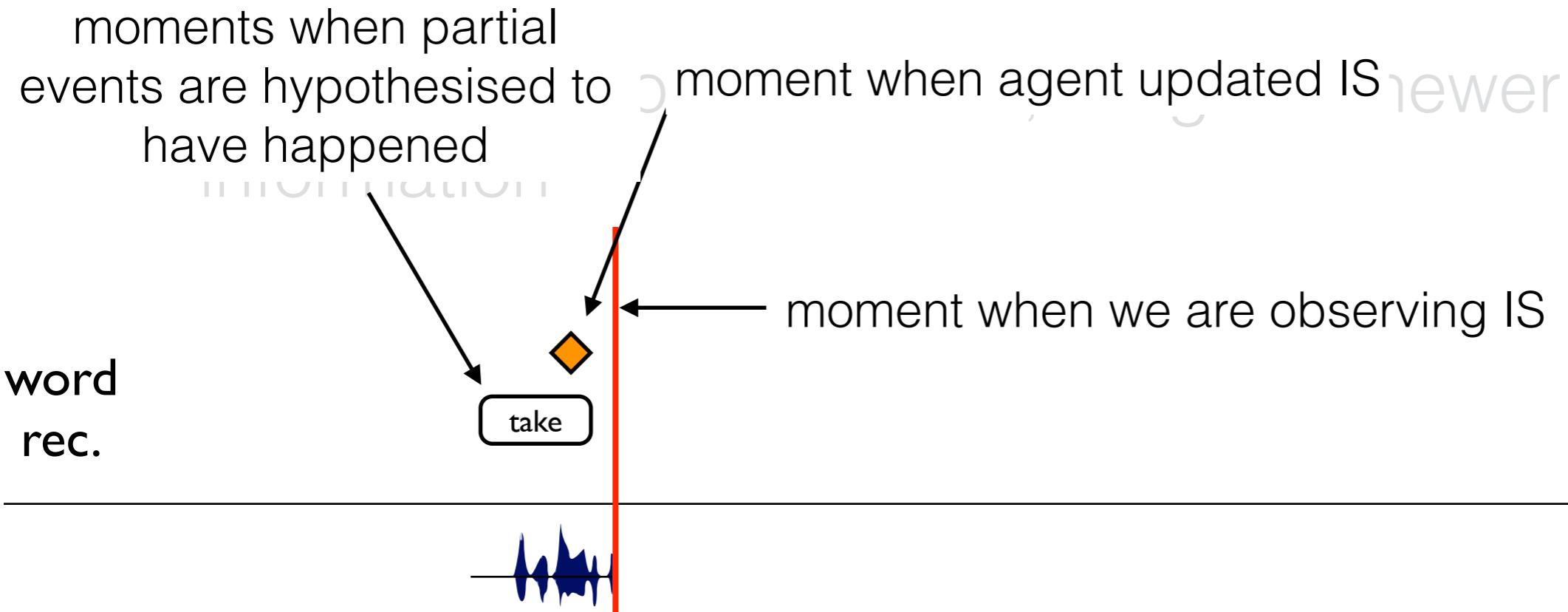
- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.
- IS may have to be revised, in light of newer information



the IU model

– Assumptions –

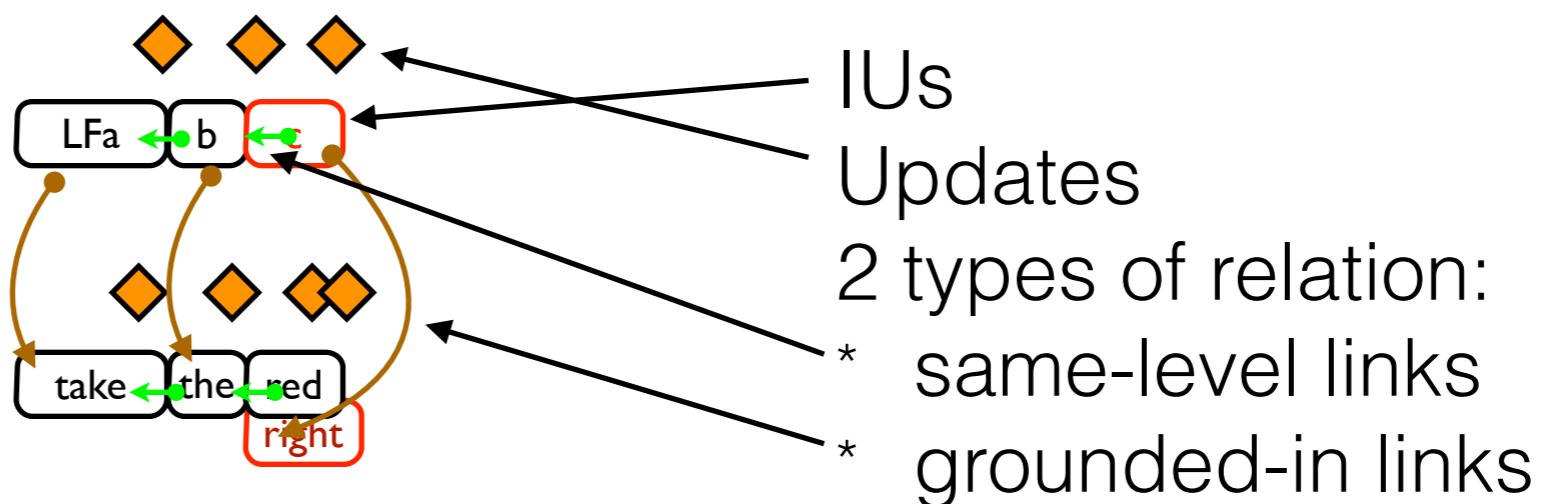
- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.



the IU model

– Assumptions –

- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.
- IS may have to be revised, in light of newer information

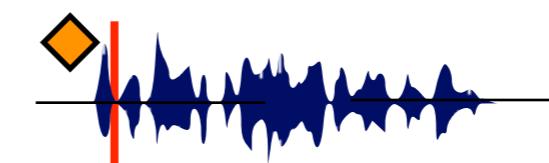


the IU model

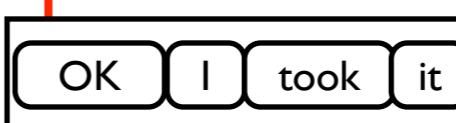
– Assumptions –

- Information state is updated with *minimal units* of information, as soon as they can be hypothesised
- “Higher-level” hypotheses can be formed on the basis of “lower-level” ones.
- IS may have to be revised, in light of newer information

TTS



NLG



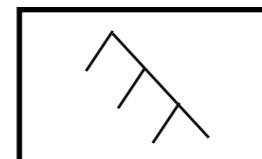
Act



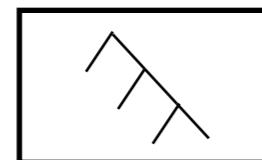
DM



Ctxt



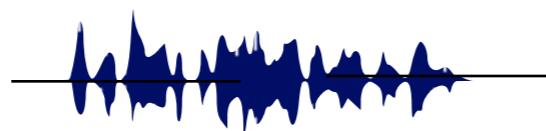
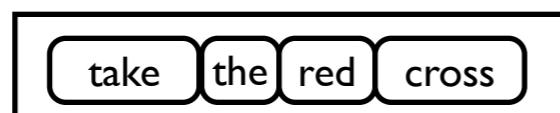
Parse

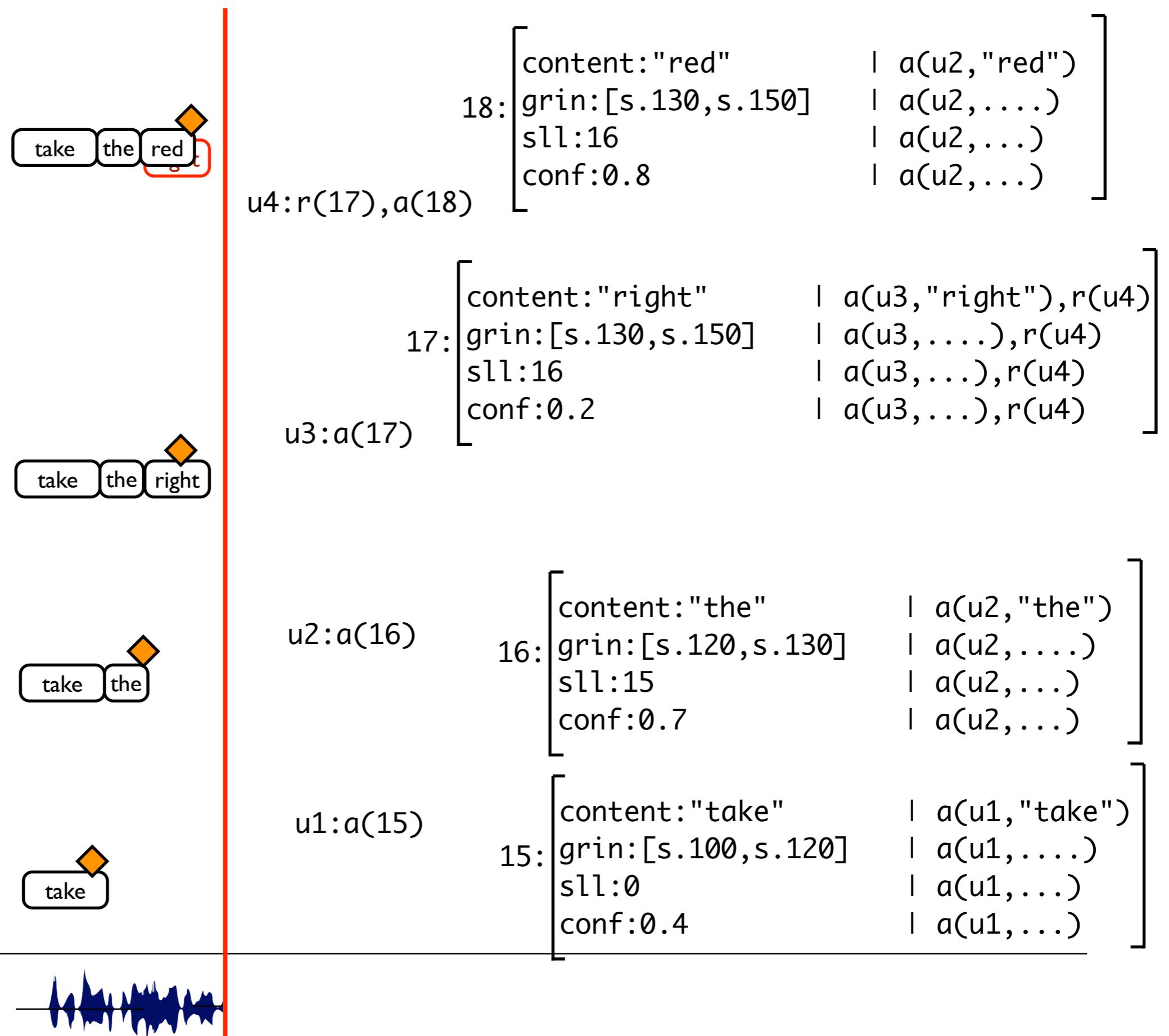


Tag



ASR

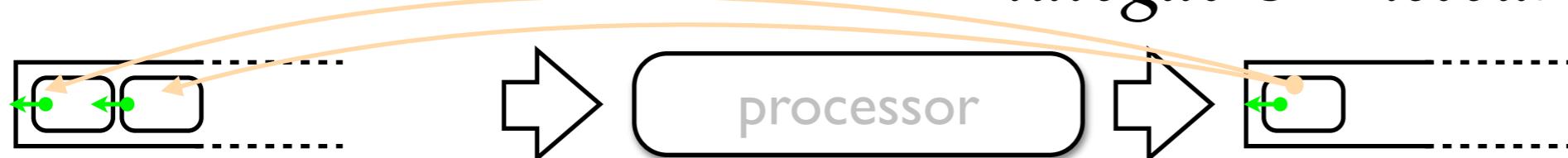




incremental processing

the IU model

- incremental units (Schlangen & Skantze, EACL 2009 /
Dialogue & Discourse 2011)



- IUs: minimal units of characteristic in/output, parts of larger unit
- part-of-relation represented as *same-level-links*
- IUs *grounded in* IUs from other levels, which were drawn on when building them
- Implemented in InproTK (<http://www.inpro.tk>), Jindigo (Skantze), IPAACA (Kopp & Buschmeier)

The NUMBERS systems fast turn-taking



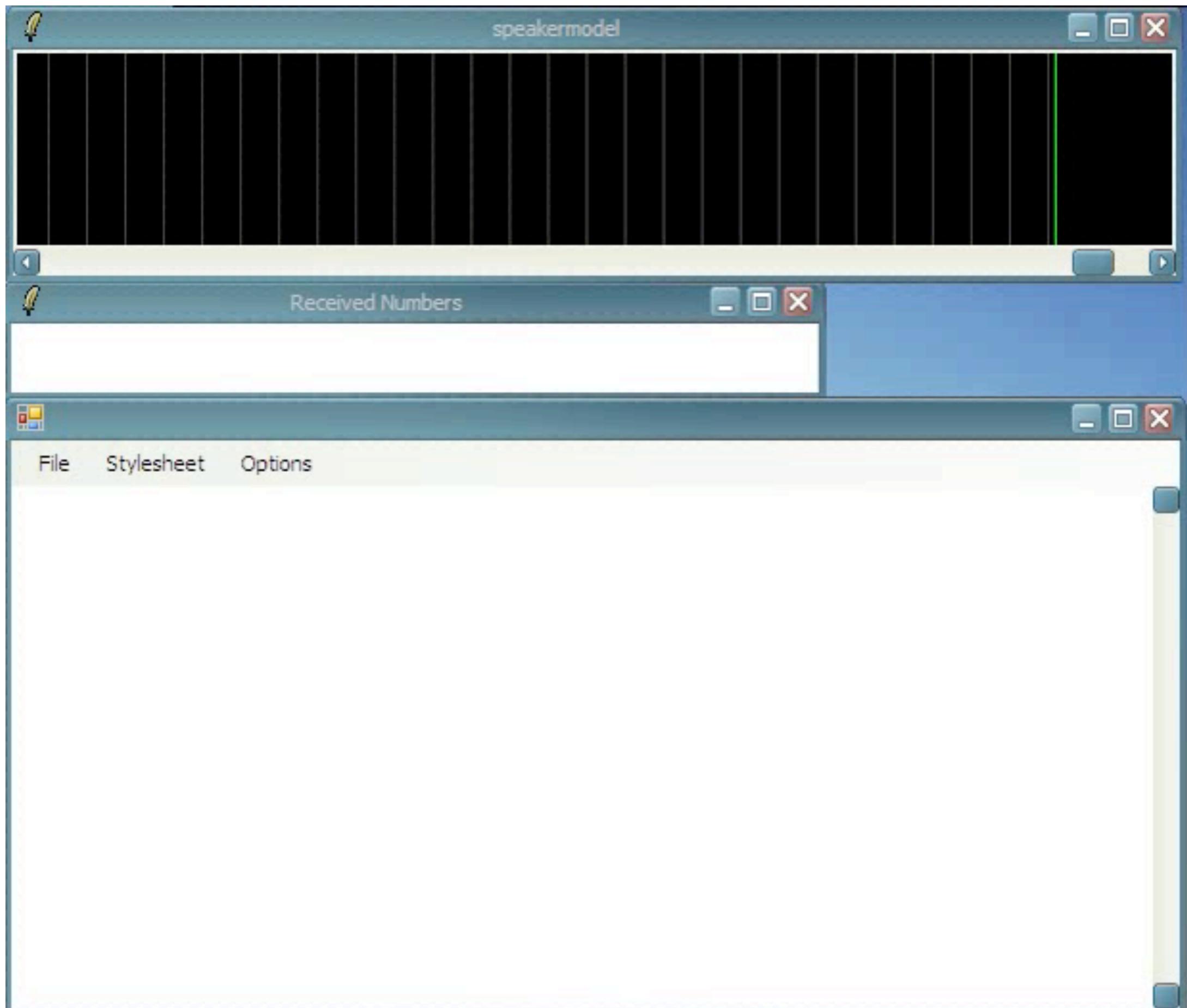
joint work with Gabriel Skantze
(Skantze & Schlangen, EACL 2009)

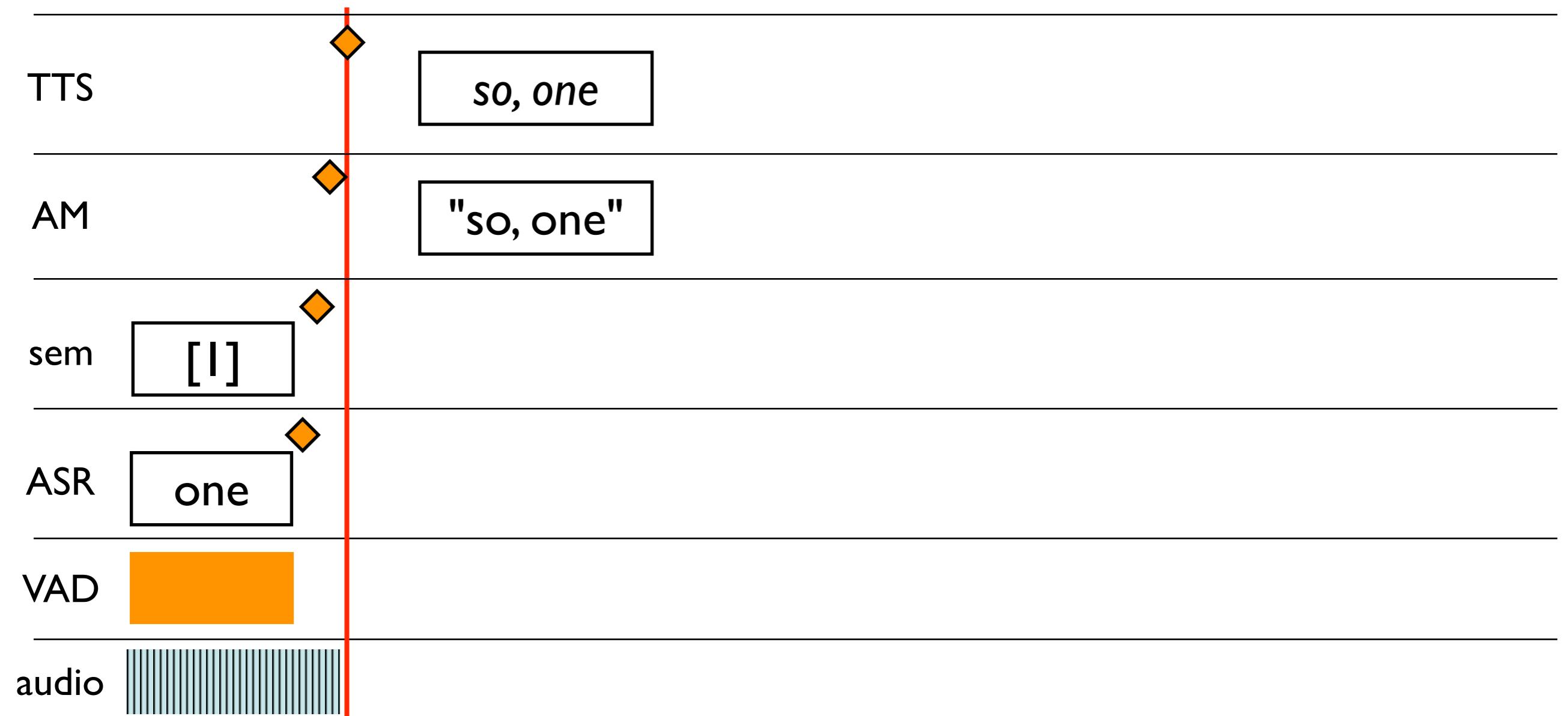
The NUMBERS systems

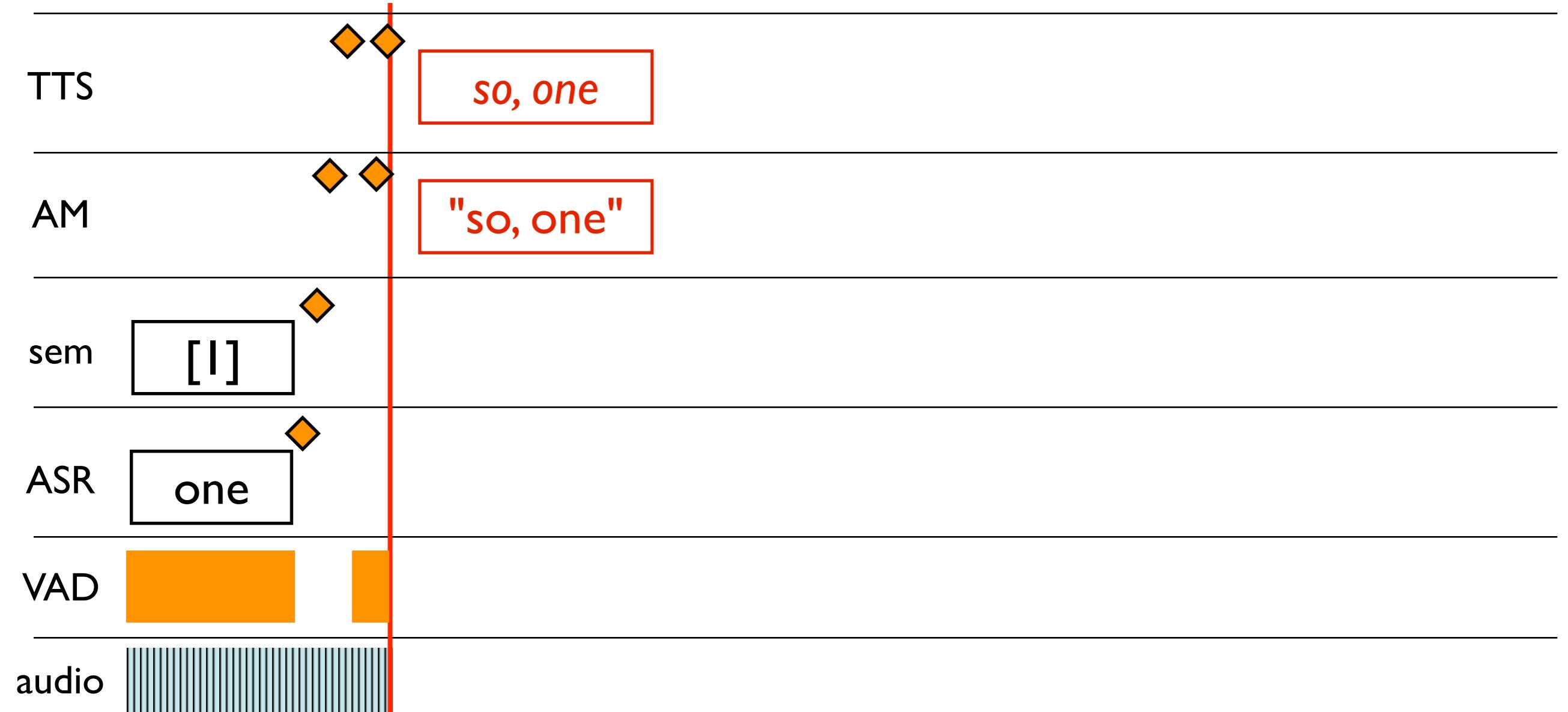
fast turn-taking

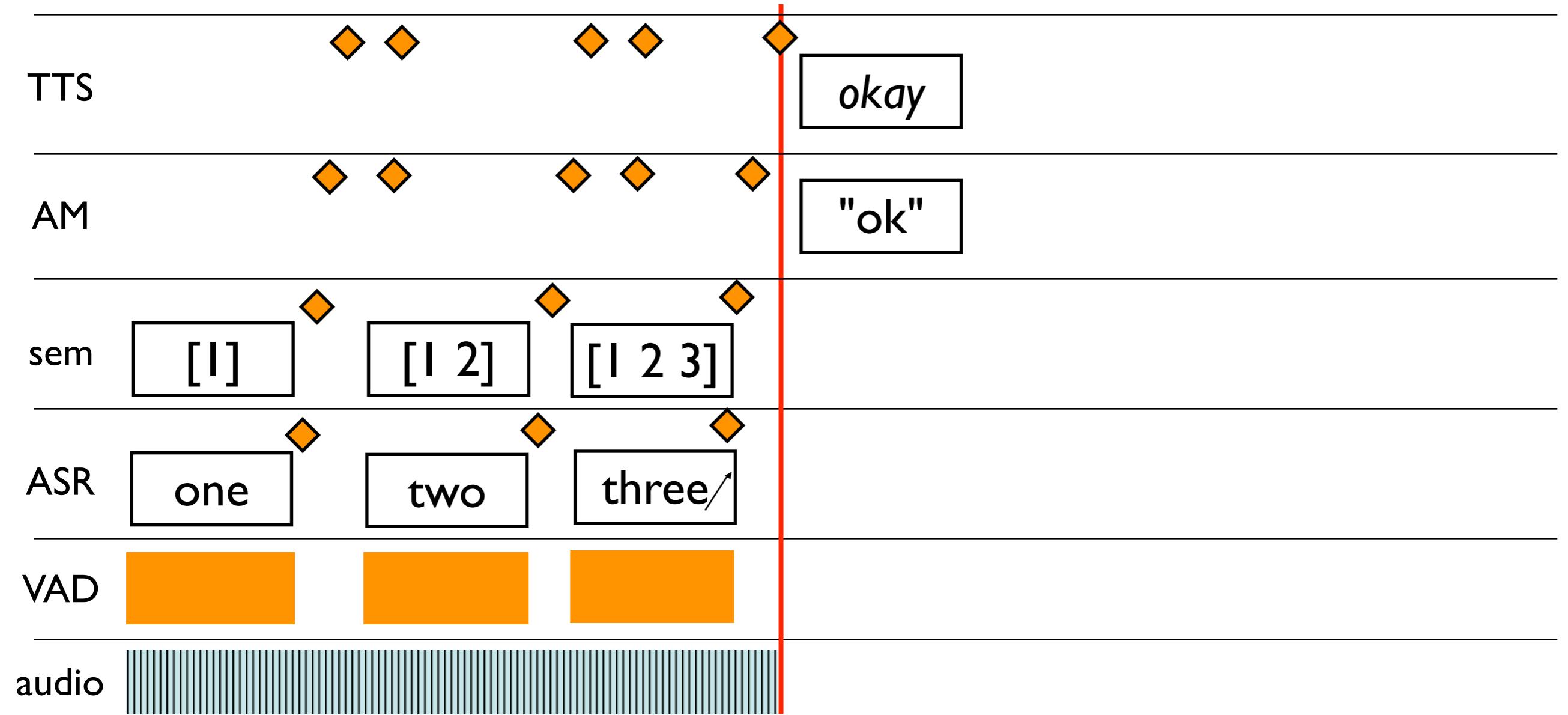
- user dictates a string of digits to system
- system tries to ground its understanding, as quickly as possible
- processing based on IU-model:
 - minimal units trigger updates
 - processors implement update functions

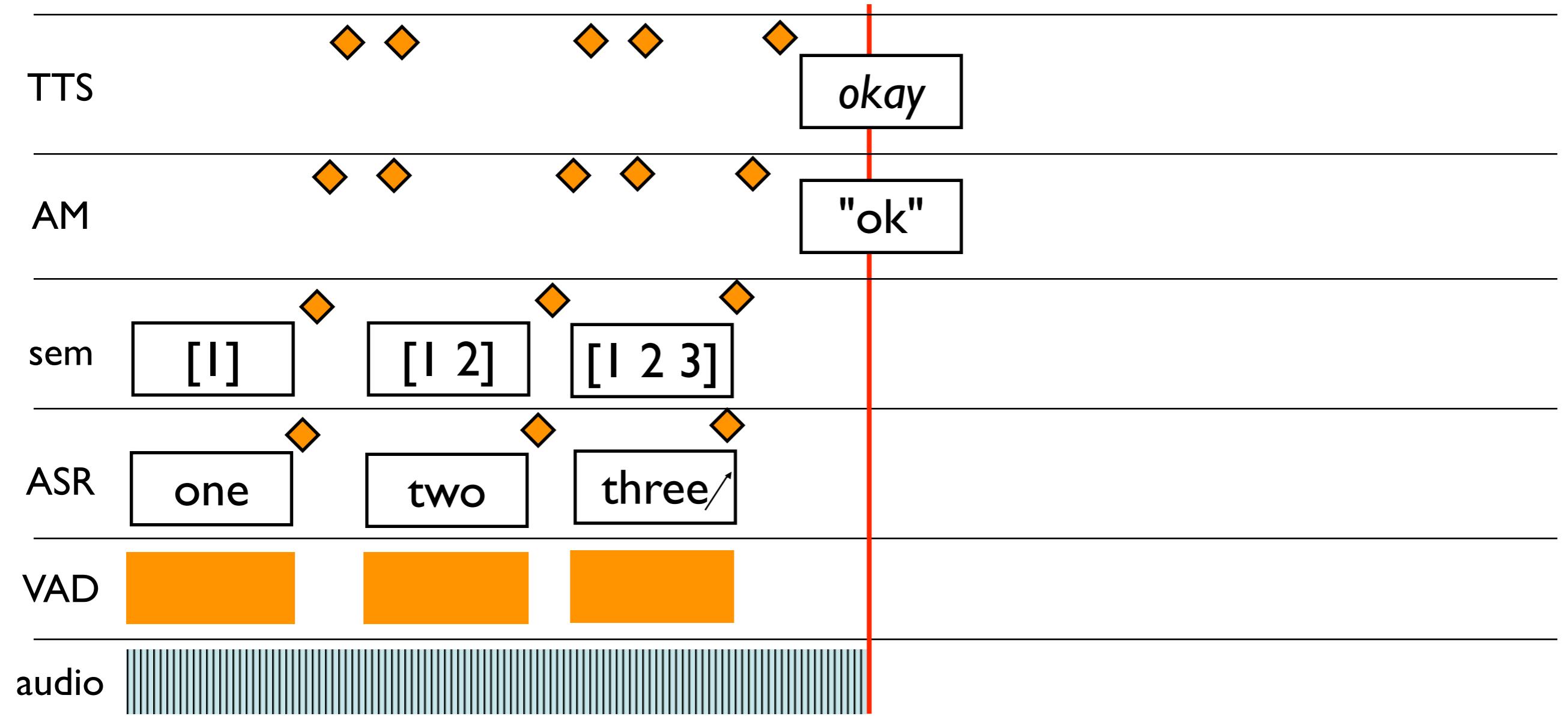
the numbers system



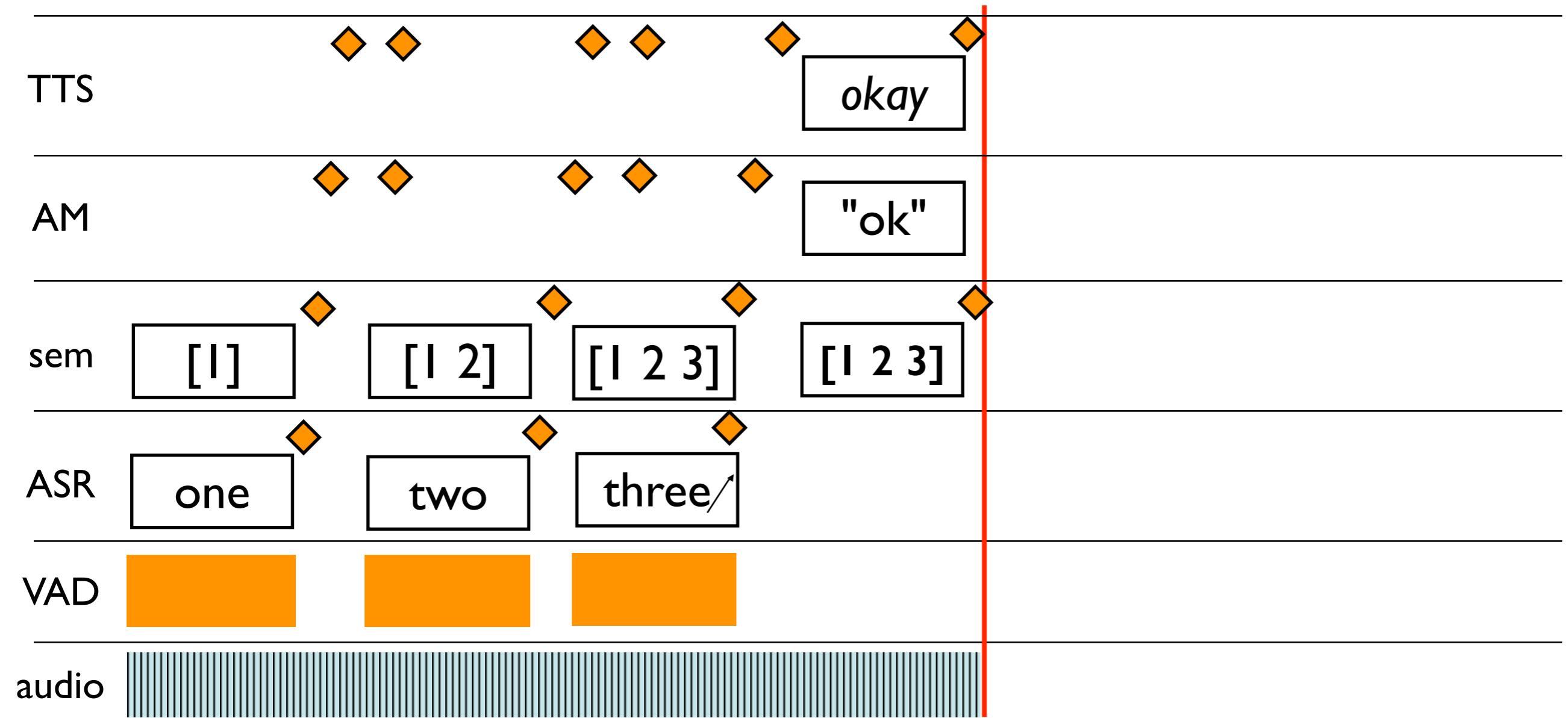








- conditional projections
- w/ dynamic offsets

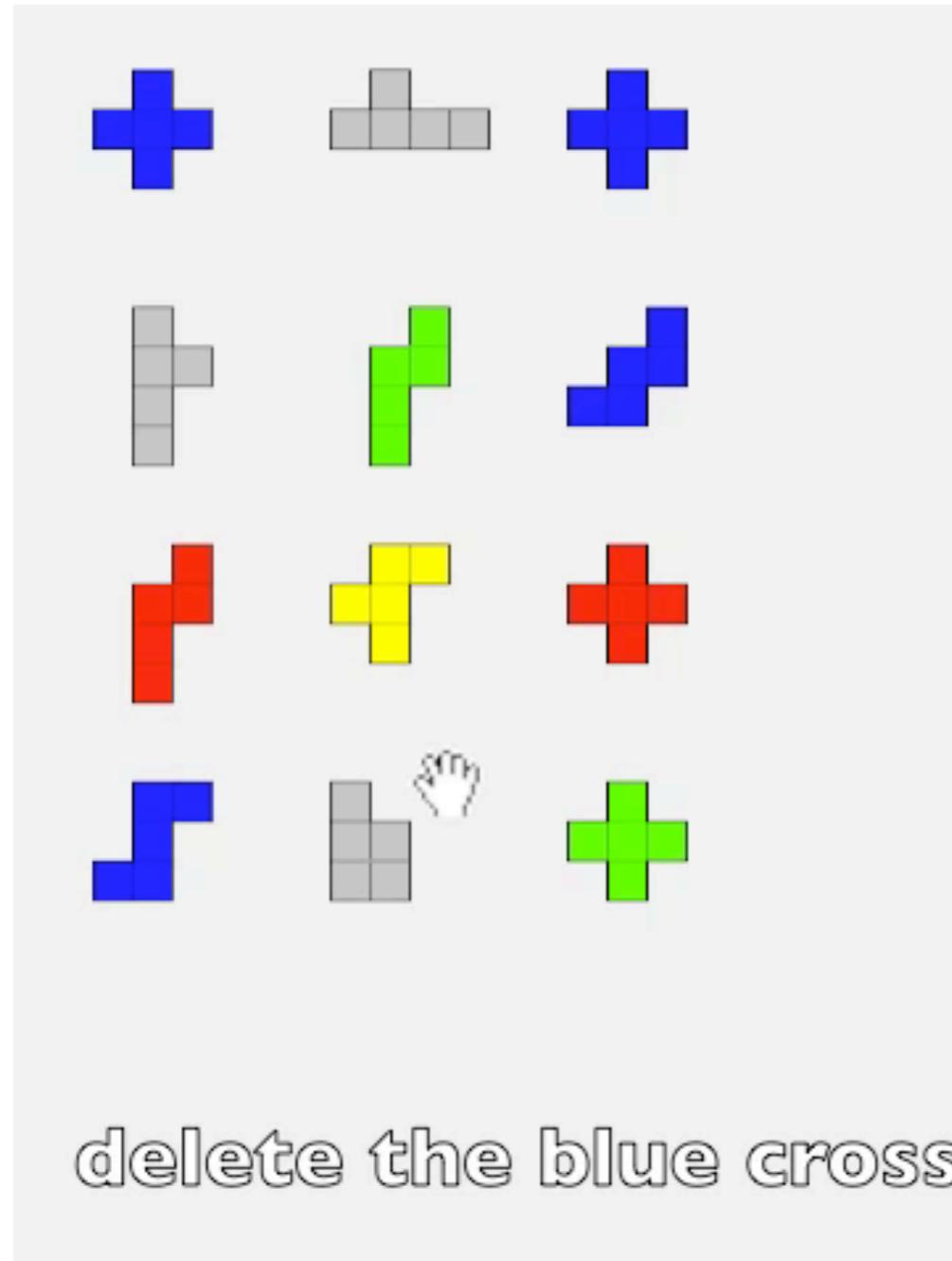


The PENTO-10 system fast turn-taking, immediate exec



joint work with Okko Buß
(Buß *et al.*, SIGdial 2010, semdial 2010, 2011)

Pentomino System



- U: *delete the blue cross*
S: which piece?
U: *top right.*
S: ok?
U: *right, now take the yellow [one]...*
S: yes?
U: *... and turn it...*
S: yes?
U: *... to the left*
S: ok.
U: *now flip the stairs...*
S: ok
U: *horizontally*
U: *that's right*
U: *erm now delete the red [one]*
S: *wh-*
U: *bottom right*
U: *correct.*

Evaluation

- Faster task completion compared to non-incremental versions of the systems
- Higher subjective ratings („would use again“, „behaves as expected“)
- Not higher task success rate

Overview

- **Part I: Foundations**

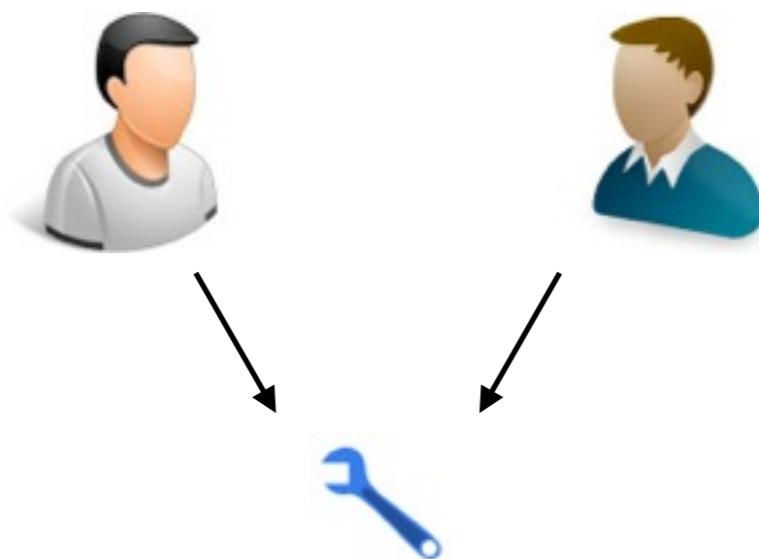
- coordination, convention
- communicative intentions
- non-conventional meaning
- grounding
- turn-taking
- disfluencies

- **Part II: Computational Models**

- approaches to dialogue modelling
- incremental processing, turn-taking
- an example: grounded semantics

overview

- the *incremental units* model of incremental dialogue processing
- realising fast turn-taking
- a simple model of incremental reference resolution

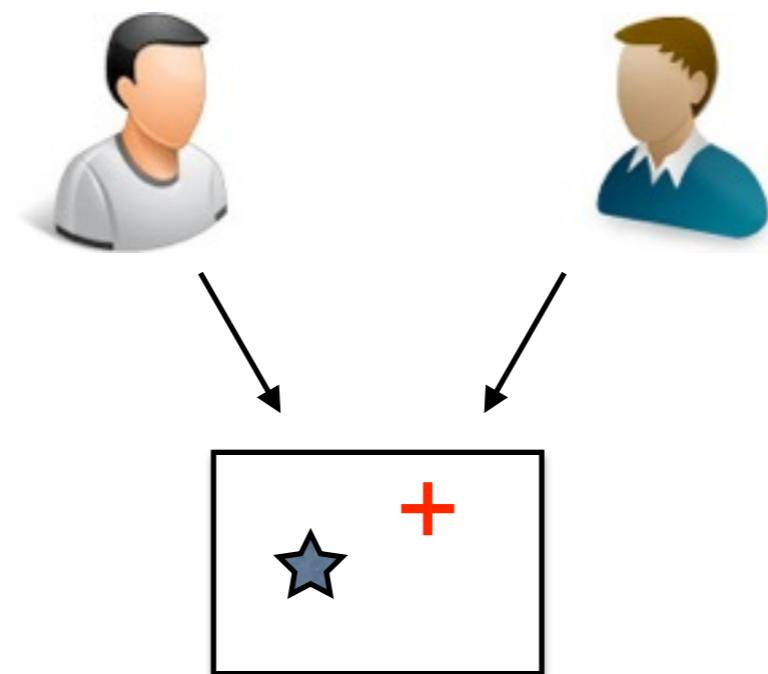


incremental statistical NLU

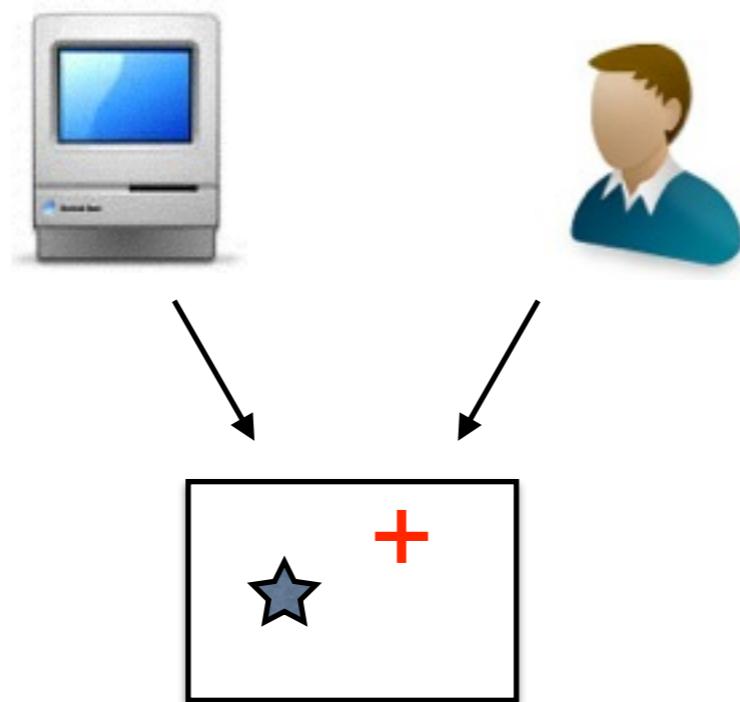


joint work with Casey Kennington
(Kennington *et al.*, SIGdial 2012, 2013; Coling 2014;
Computer Speech & Language 2014)

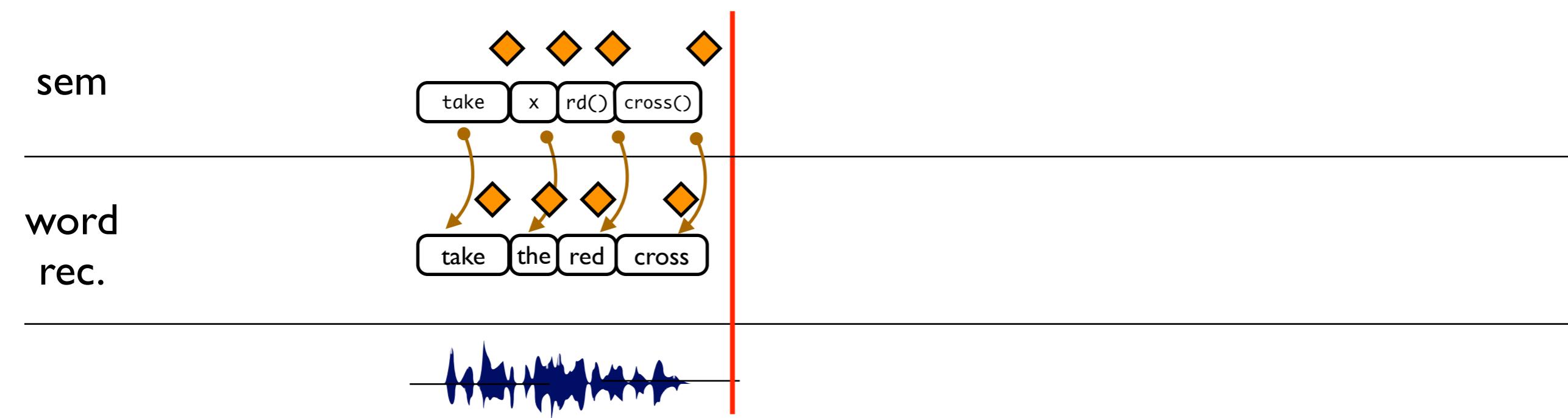
incremental statistical NLU



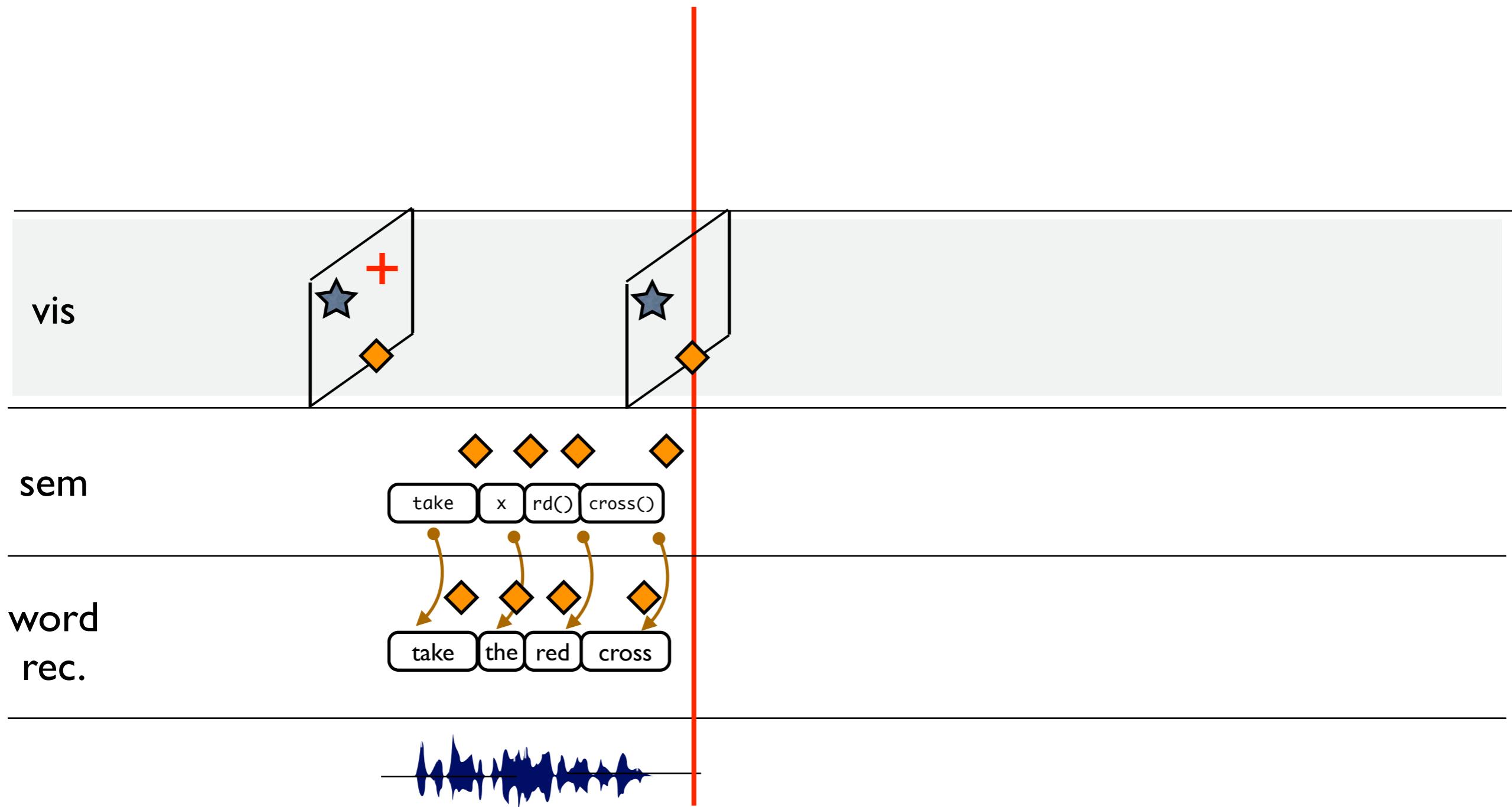
incremental statistical NLU



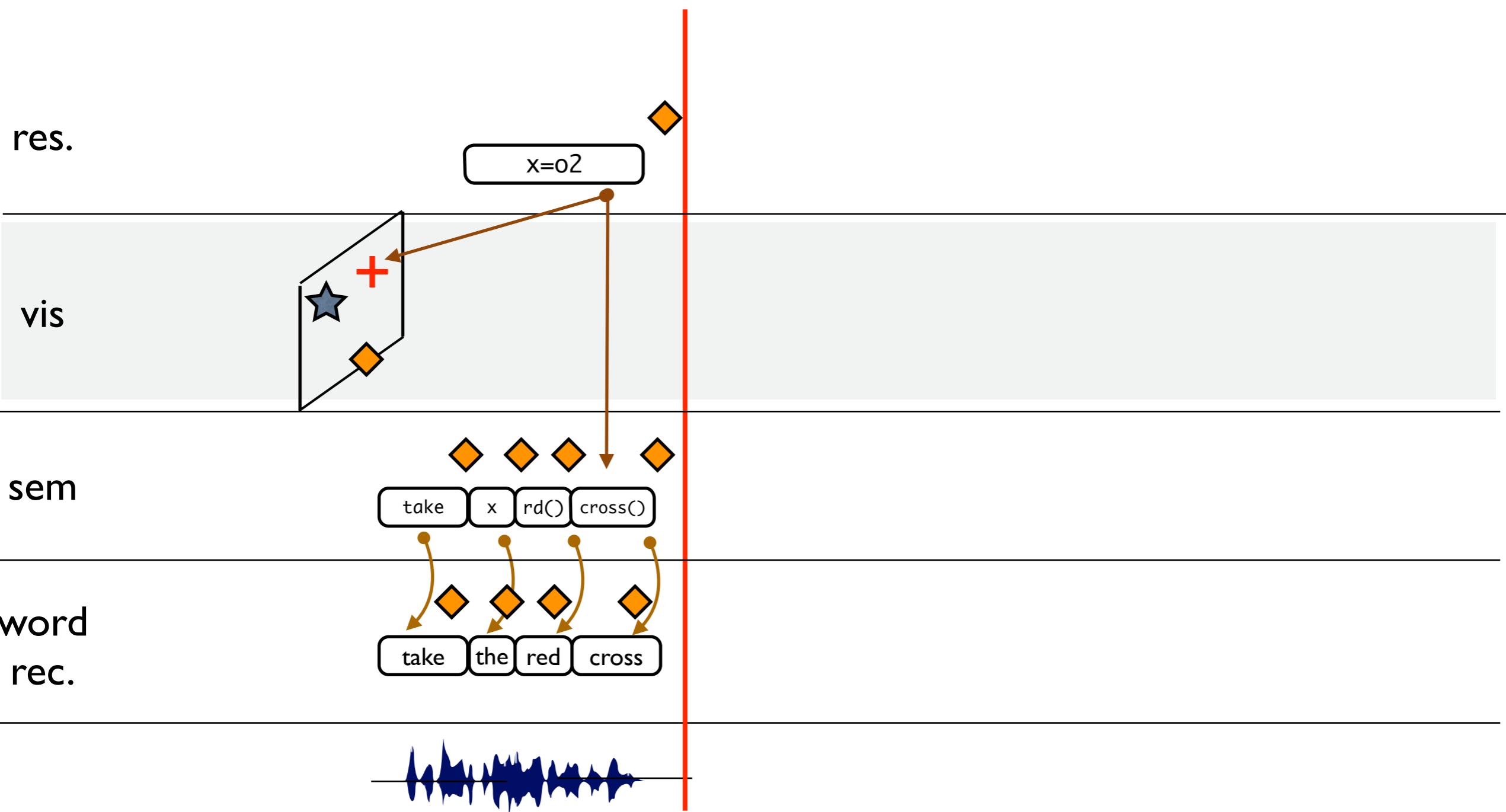
incremental statistical NLU



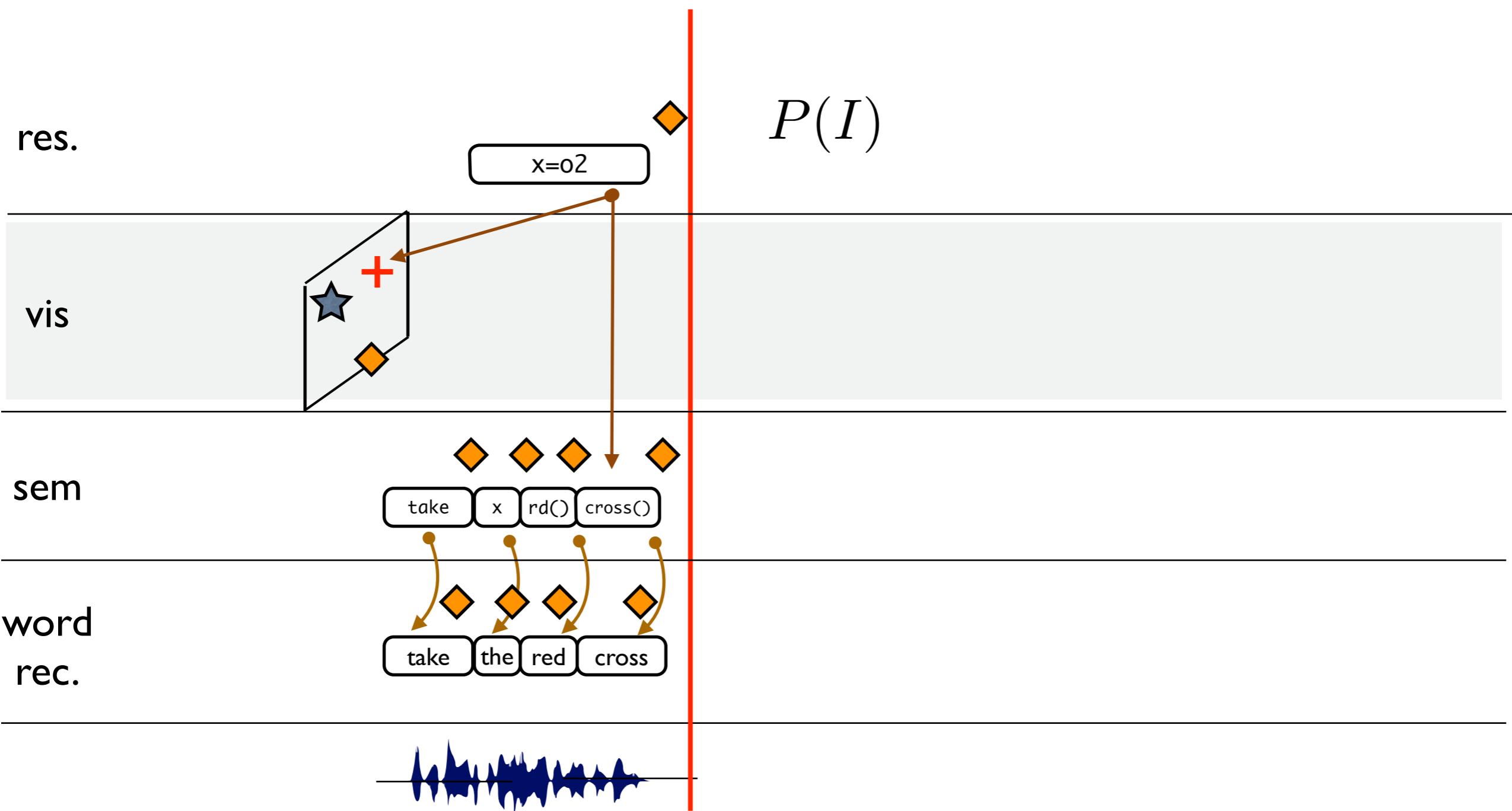
incremental statistical NLU



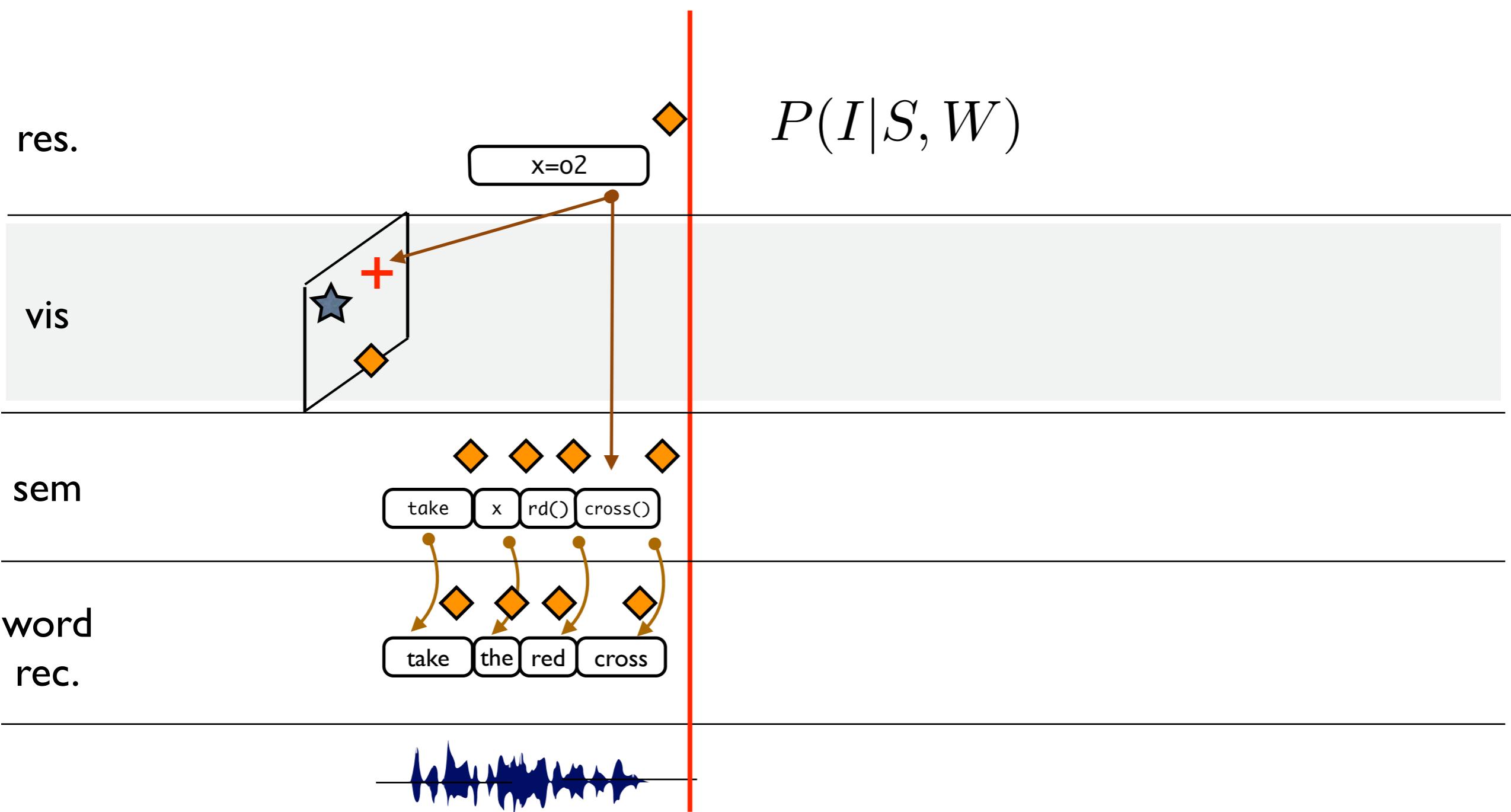
incremental statistical NLU



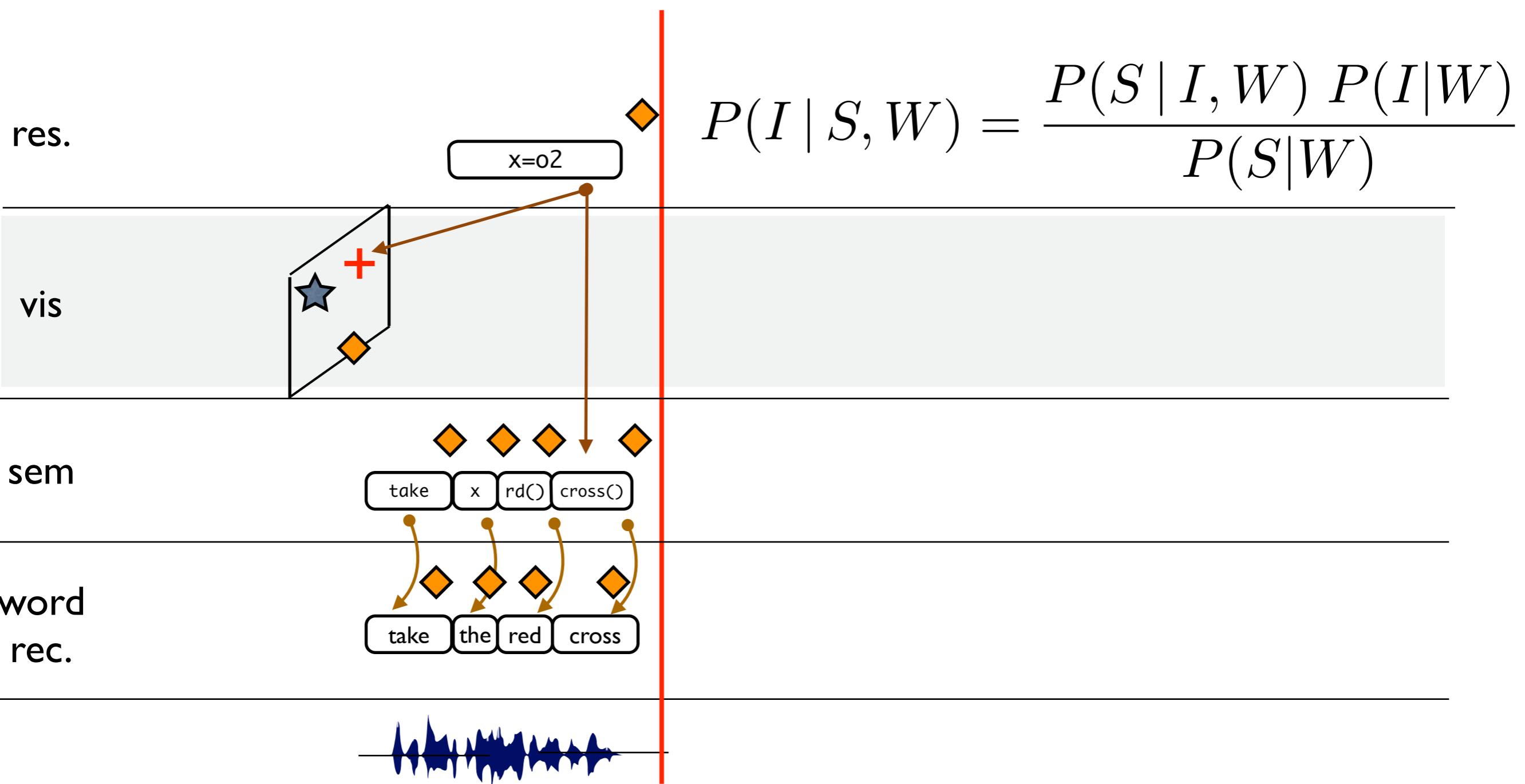
incremental statistical NLU



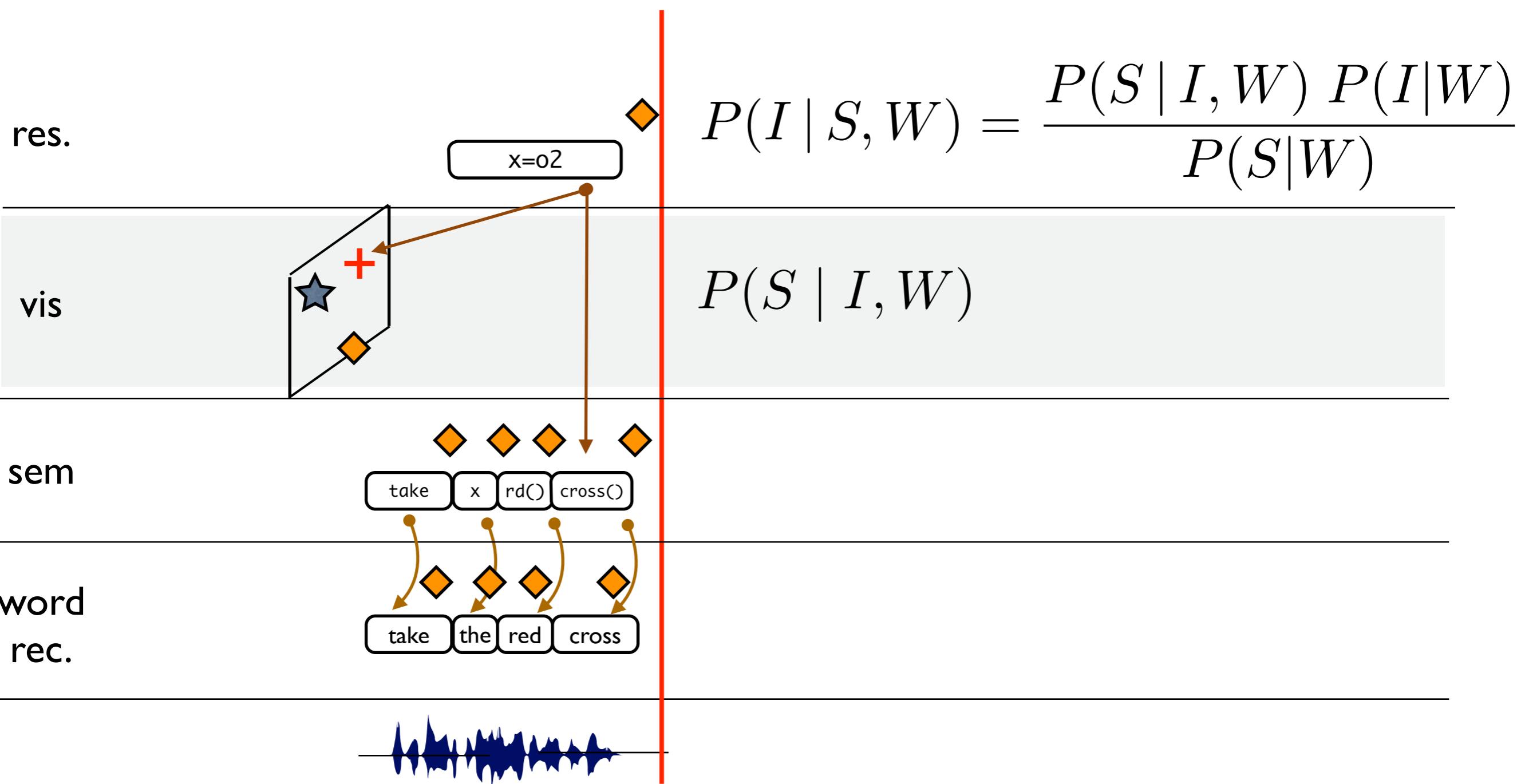
incremental statistical NLU



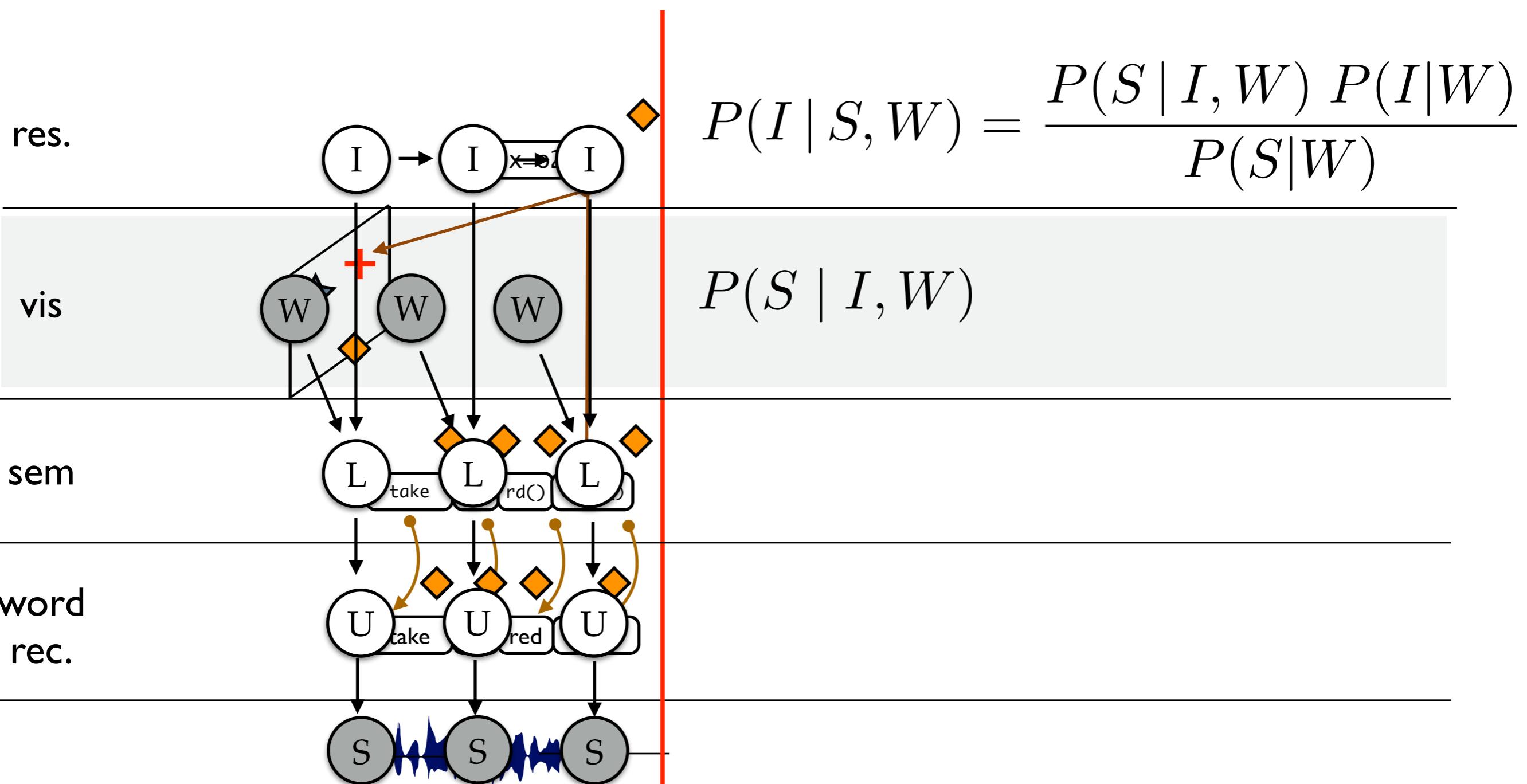
incremental statistical NLU



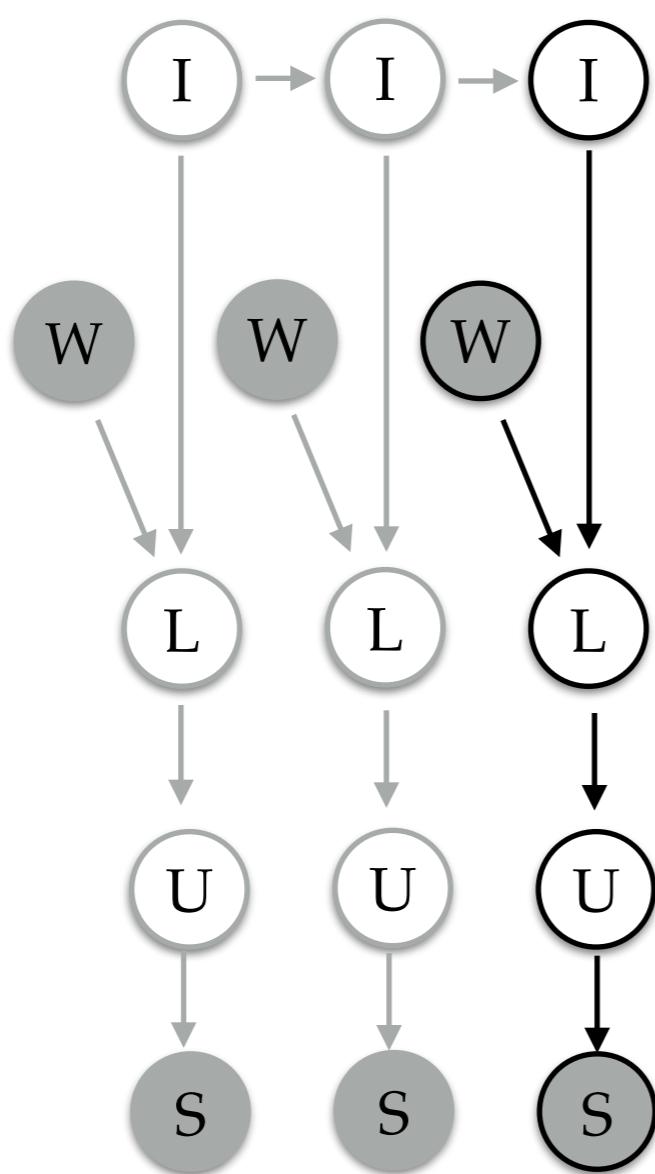
incremental statistical NLU



incremental statistical NLU

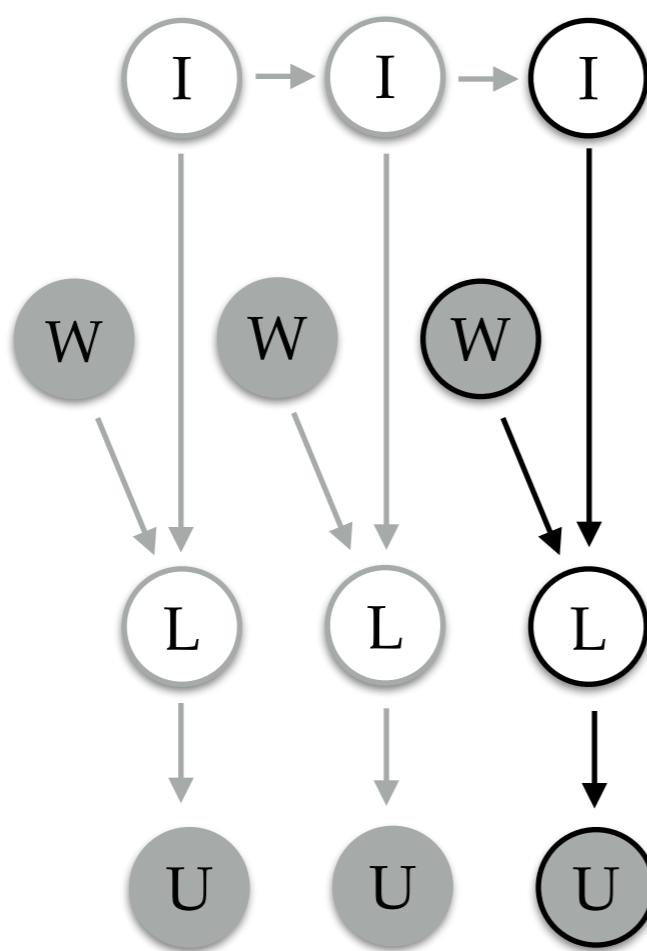


incremental statistical NLU



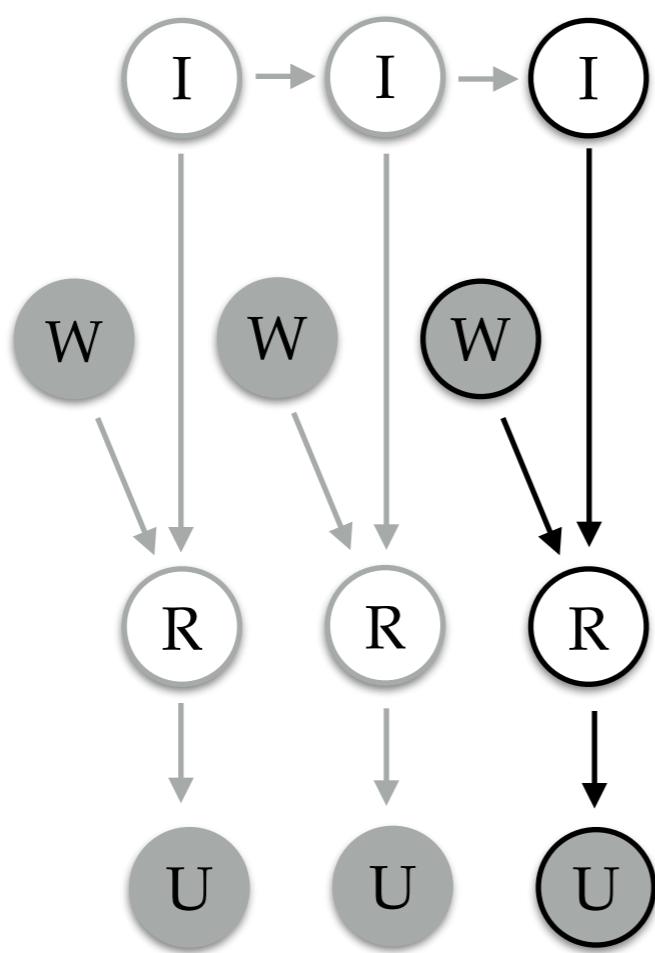
$$P(I | S, W) = \frac{P(S | I, W) P(I|W)}{P(S|W)}$$
$$P(S | I, W)$$

incremental statistical NLU



$$P(I | S, W) = \frac{P(S | I, W) P(I|W)}{P(S|W)}$$
$$P(U | I, W)$$

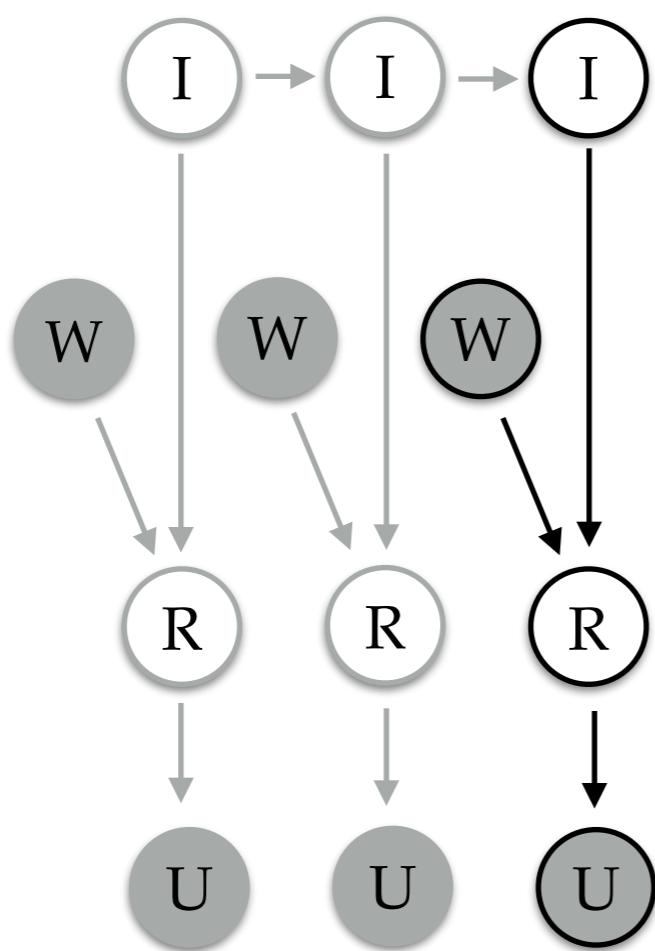
incremental statistical NLU



$$P(I | S, W) = \frac{P(S | I, W) P(I|W)}{P(S|W)}$$

$$P(U|I, W) = \sum_R P(U, R|I, W)$$

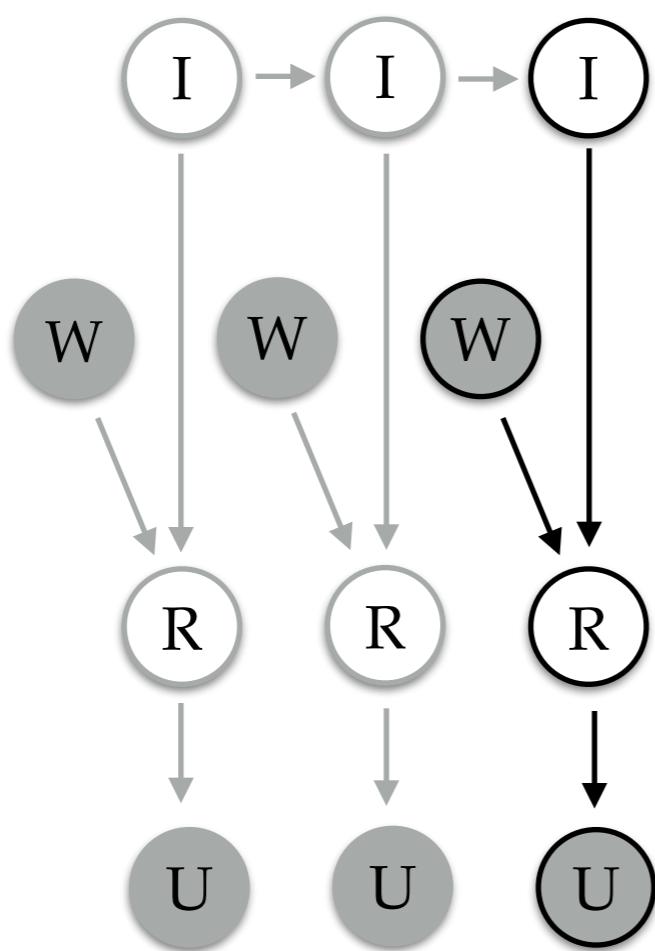
incremental statistical NLU



$$P(I | S, W) = \frac{P(S | I, W) P(I|W)}{P(S|W)}$$

$$\begin{aligned} P(U|I, W) &= \\ \sum_R P(U, R|I, W) &= \\ \sum_R P(U|R)P(R|I, W) \end{aligned}$$

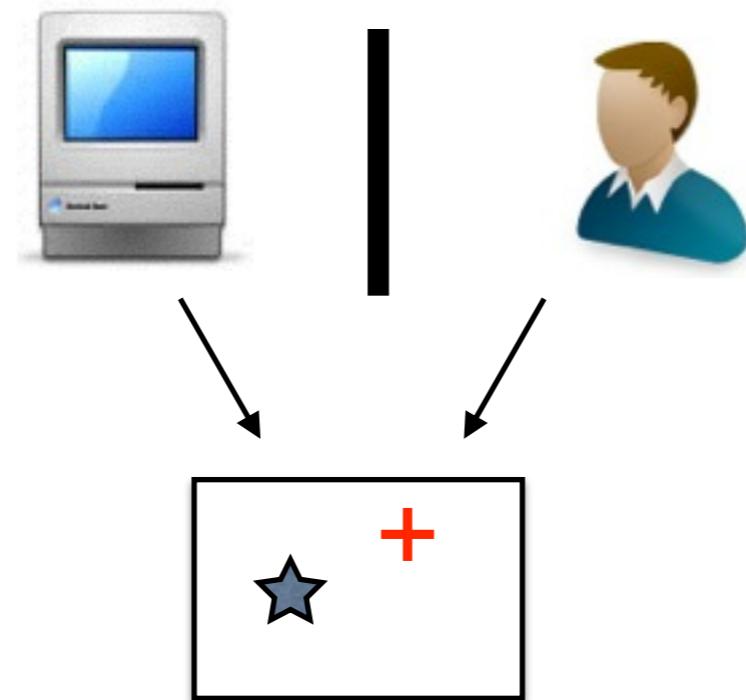
incremental statistical NLU



$$\begin{aligned} P(I|U, W) &= \\ P(U|I, W)P(I|W) \frac{1}{P(U|W)} &= \\ \sum_R P(U|R)P(R|I, W) \\ P(I|W) \frac{1}{P(U|W)} \end{aligned}$$

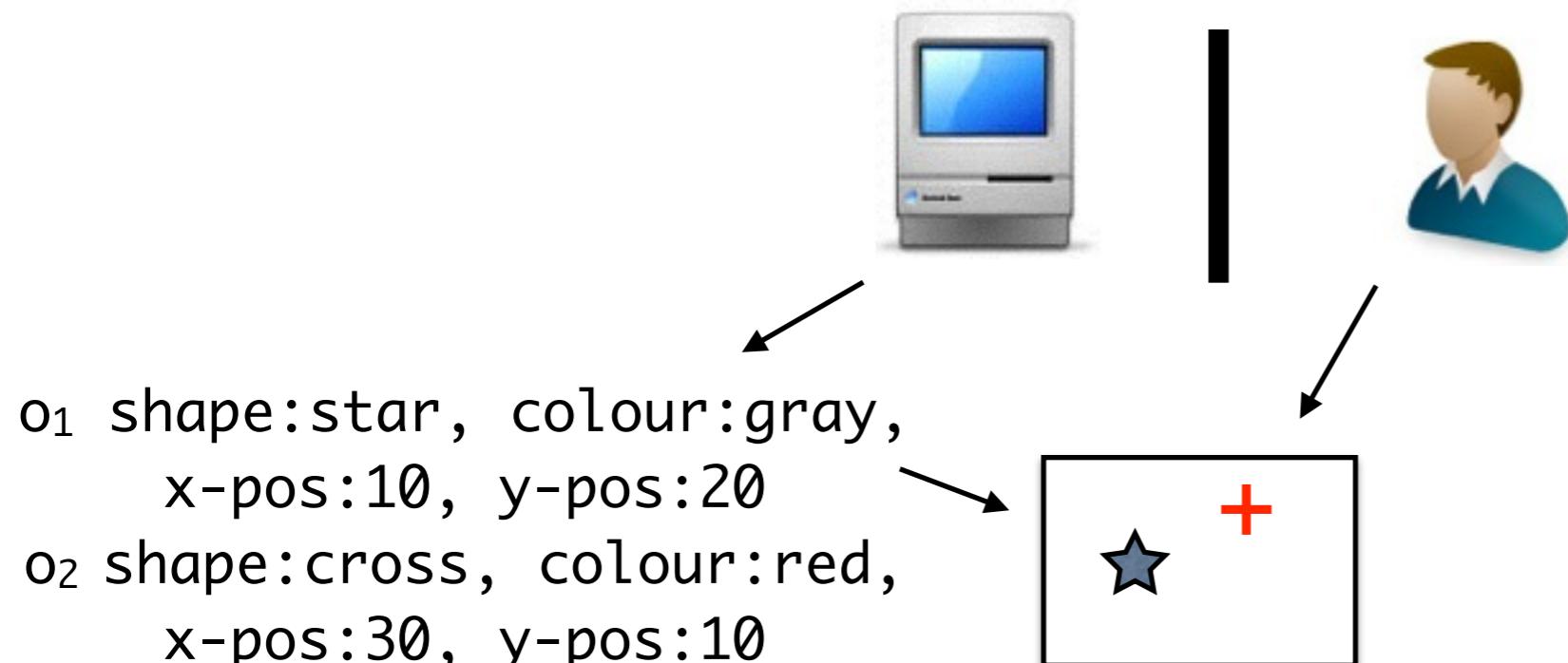
incremental statistical NLU

– The Pento-2010 Data –



incremental statistical NLU

– The Pento-2010 Data –

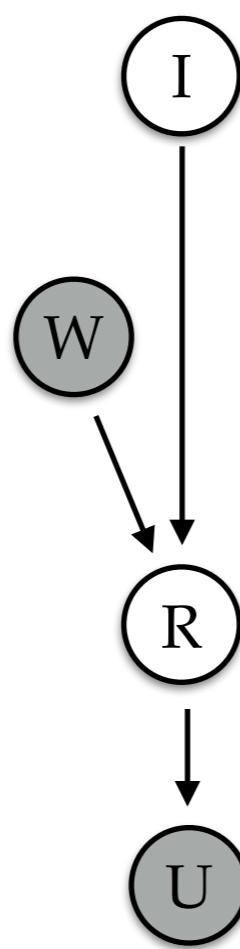


	o	o
star	1	0
cross	0	1
gray	1	0
red	0	1

incremental statistical NLU

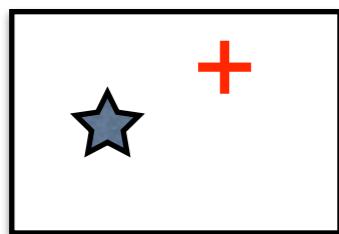
	o	o
star	1	0
cross	0	1
gray	1	0
red	0	1

	o	o
star	.5	0
cross	0	.5
gray	.5	0
red	0	.5

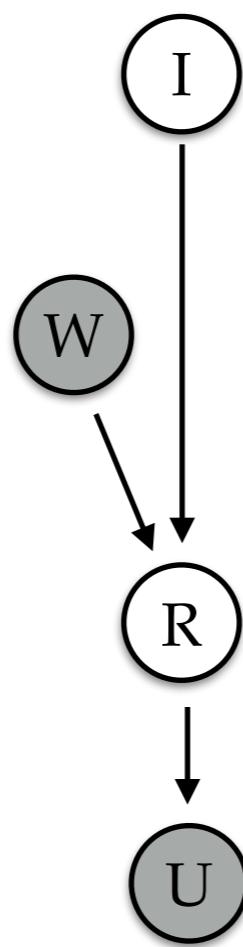


$$P(I|U, W) = \frac{1}{P(U|W)} P(I|W) \sum_R P(R|I, W) P(U|R)$$

incremental statistical NLU



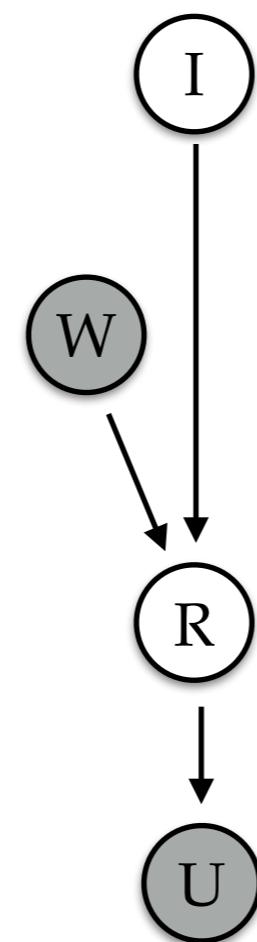
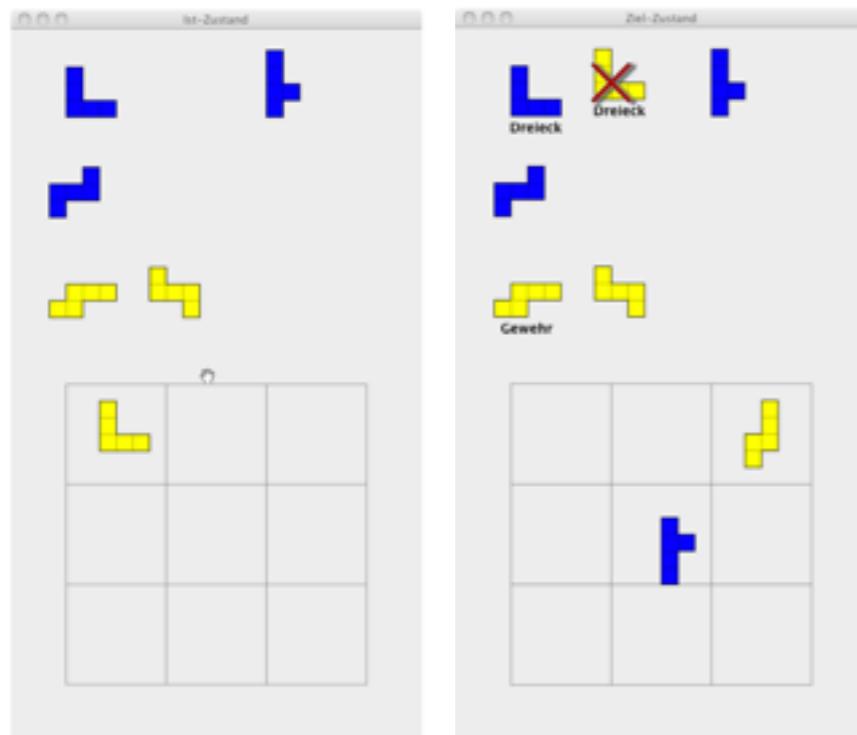
	o	o
star	.5	0
cross	0	.5
gray	.5	0
red	0	.5



$$P(I = o_2 | U = \text{cross}, W) = \frac{1}{P(U=\text{cross}|W)} P(I = o_2 | W)$$
$$\sum_R P(R|I = o_2, W) P(U = \text{cross}|R)$$

incremental statistical NLU

– Results for Pento-2010 Data –

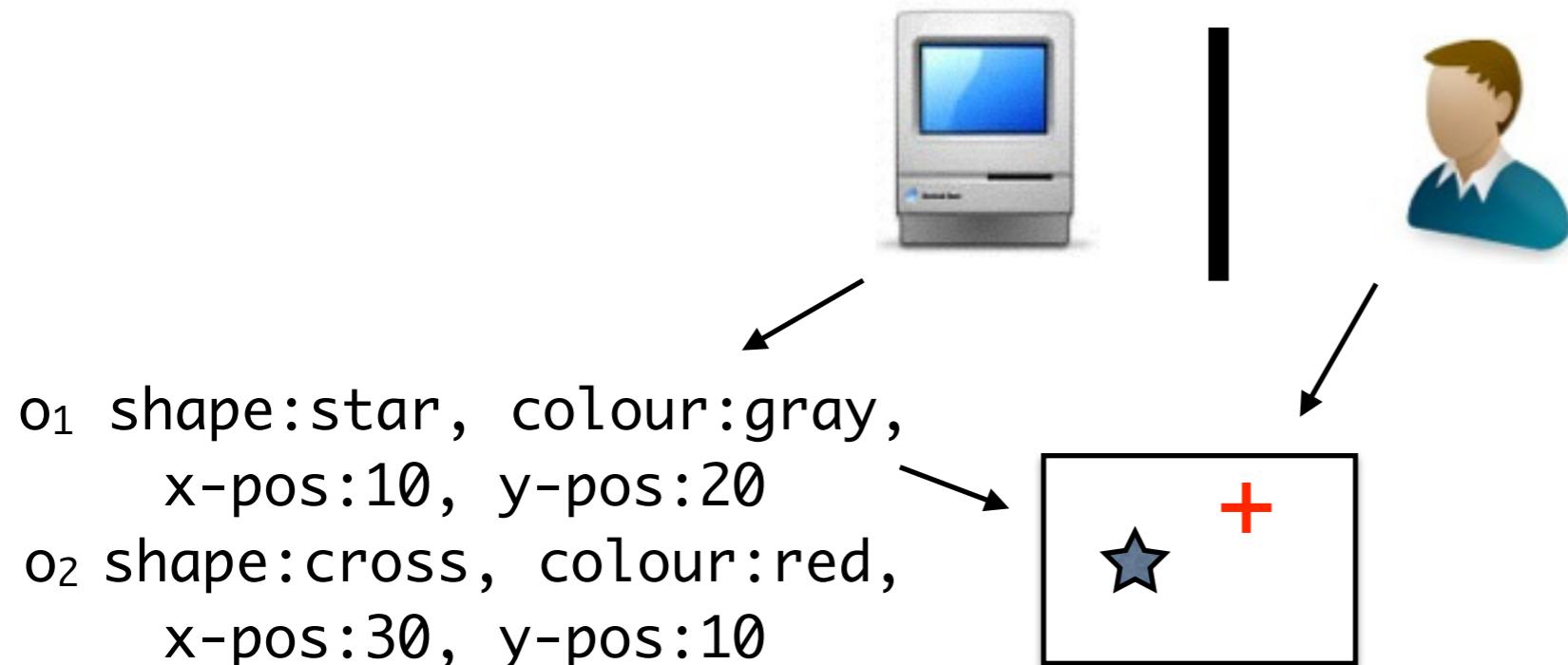


$$P(I|U, W) = \frac{1}{P(U|W)} P(I|W) \sum_R P(R|I, W) P(U|R)$$

	NB	ME	K	H	P
fscore	81.16 (74.5)	92.26 (89.4)	92.18 (86.8)	76.9	
slot	73.62 (66.4)	88.91 (85.1)	88.88 (81.6)		
frame	42.57 (34.2)	74.08 (67.2)	74.76 (61.2)		
action	80.05	93.62	92.62		
object	76.22	90.79	84.71		
result	64.4	82.54	86.65	64.3	

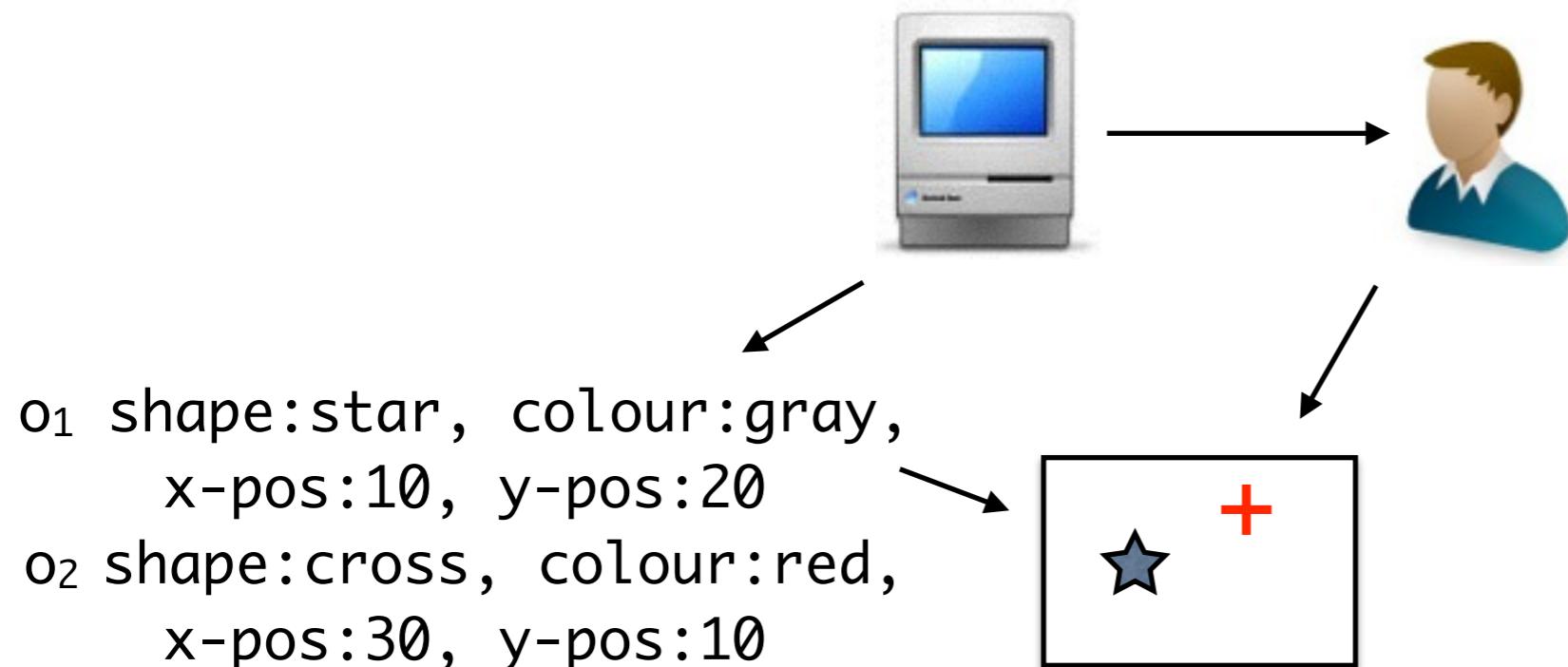
incremental statistical NLU

– The Pento-2010 Data –



incremental statistical NLU

– The Pento-2013 Data –



multimodal processing

– mint.tools / venice –

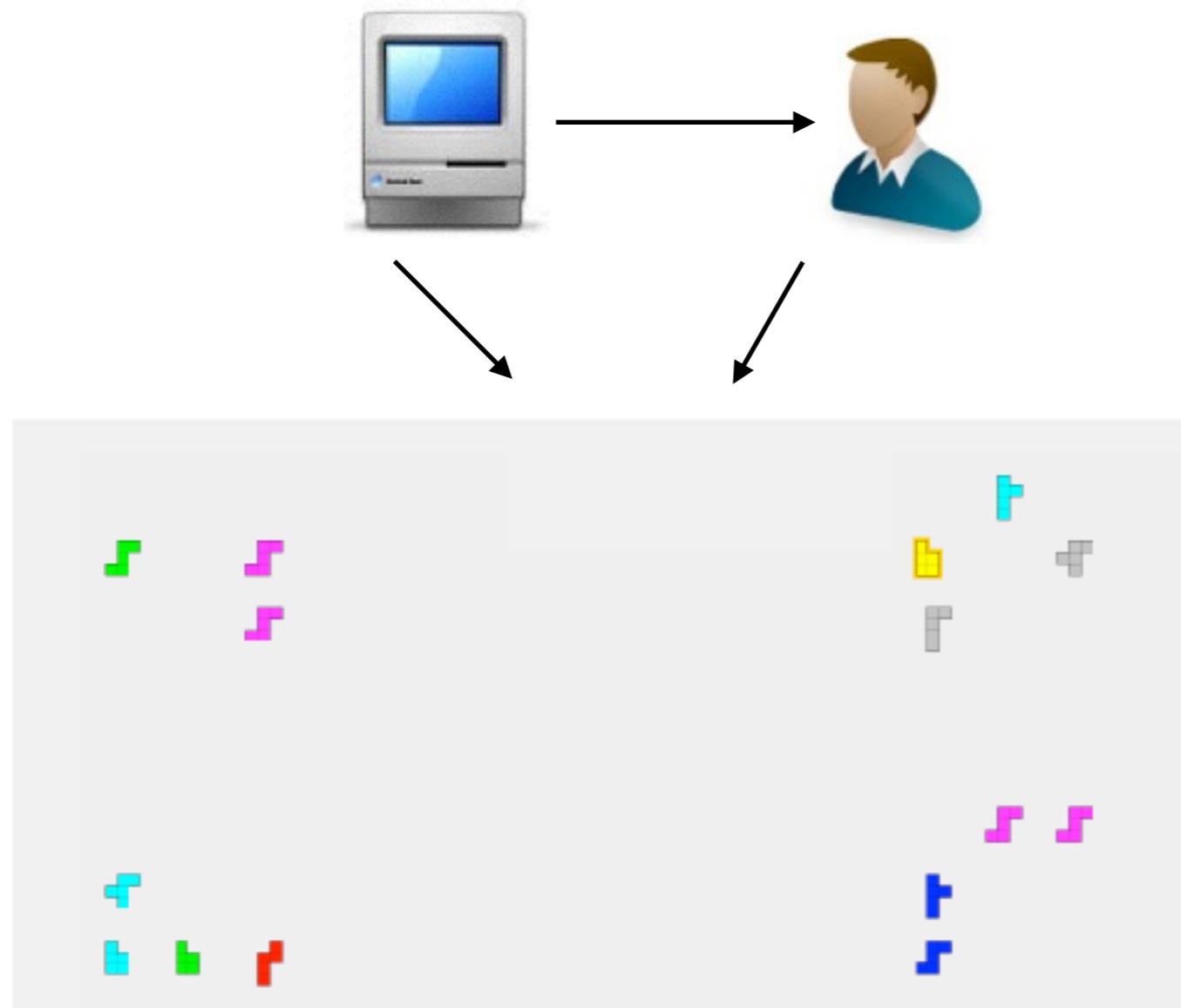


joint work with Spyros Kousidis
(Kousidis *et al.*, SIGdial 2013, 2014; ICMI 2014)

technical infrastructure for combining various sensors (motion capture, eye tracking) and various actuators (NAO robot) with InproTK

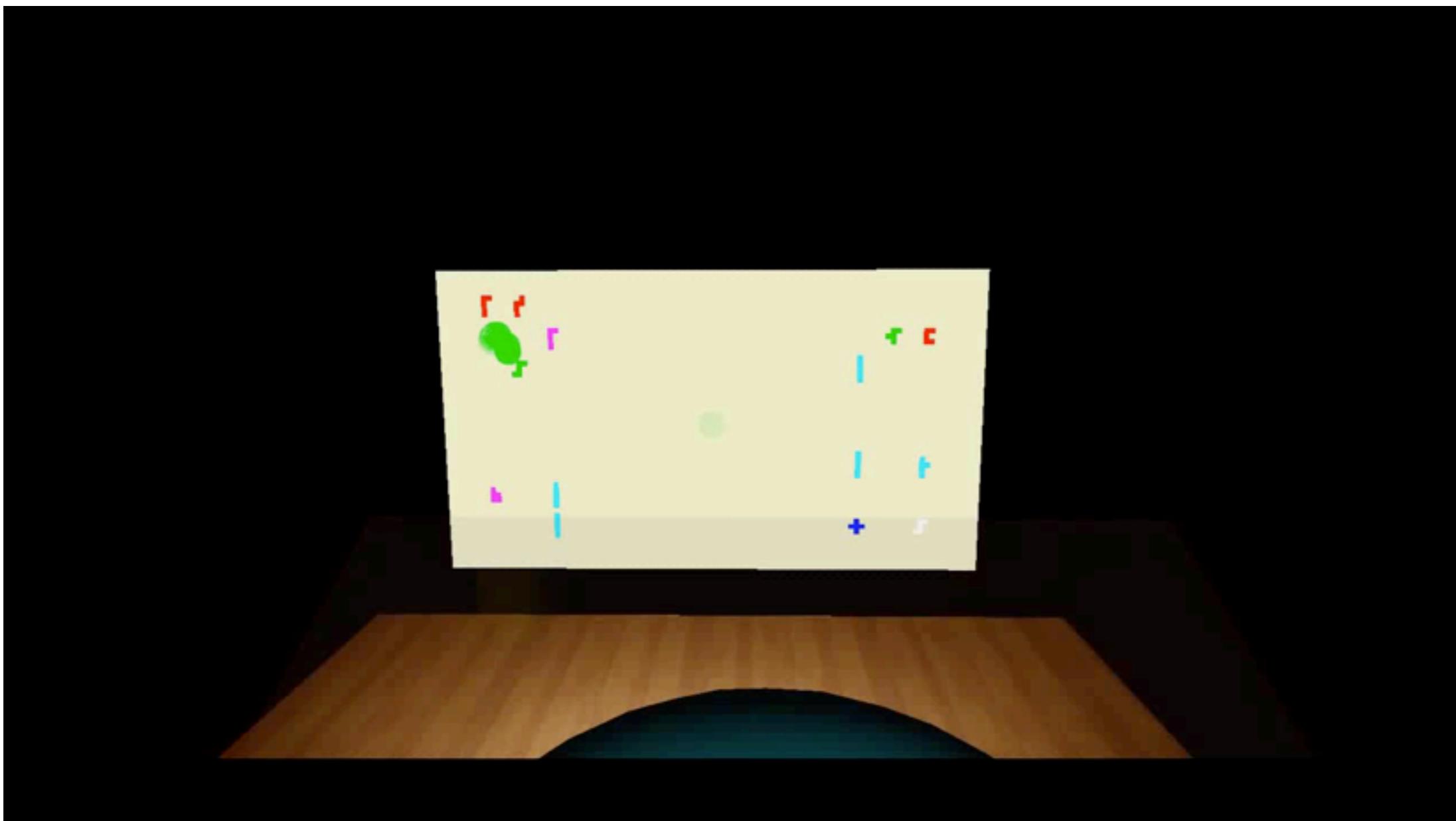
incremental statistical NLU

– The Pento-2013 Data –



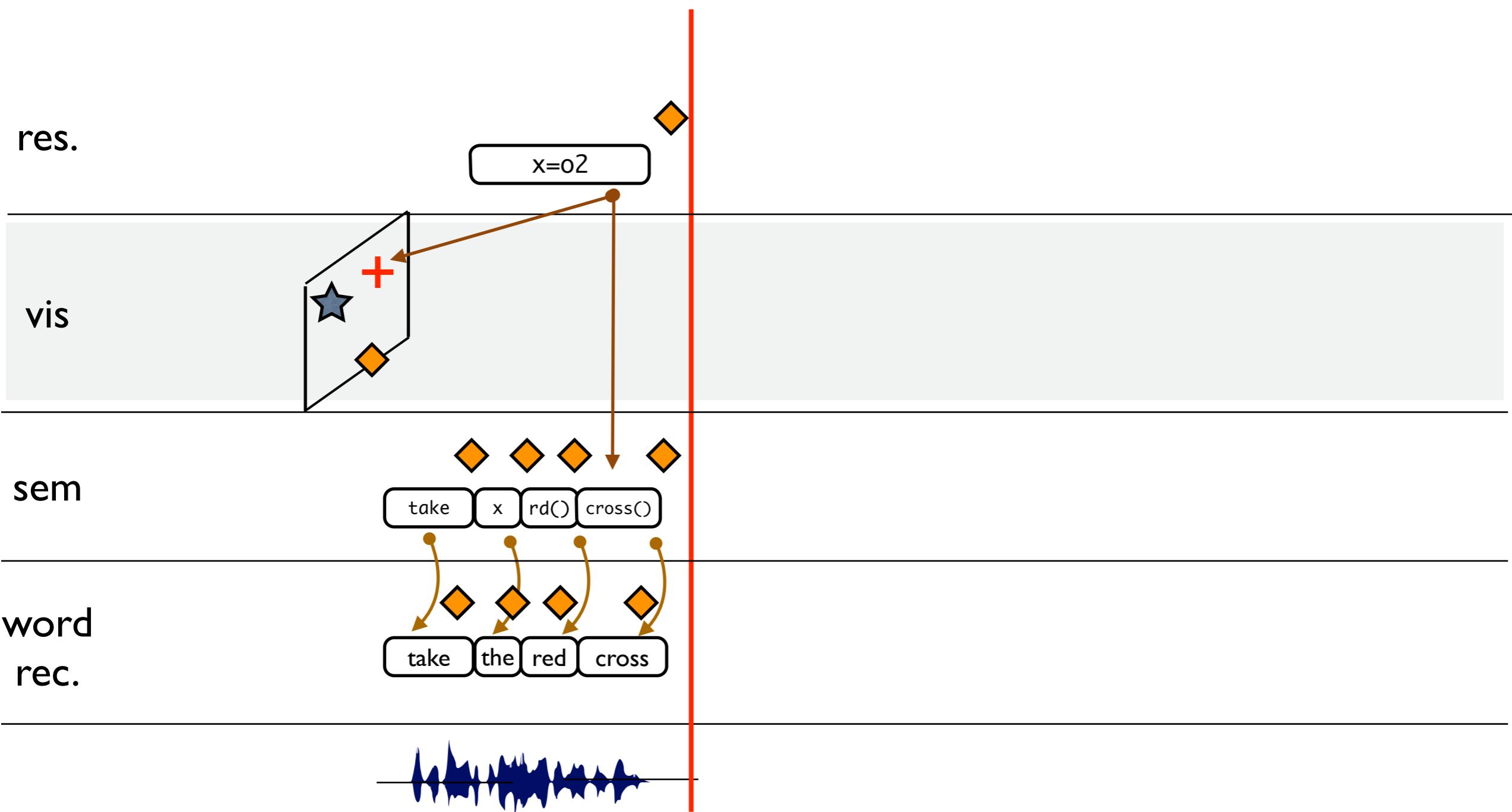
incremental statistical NLU

– The Pento-2013 Data –



incremental statistical NLU

– Adding Gaze & Gestures –



Properties over Time

dann ... nehmen ... zweite t bisschen rechts ist ... rüssel hobel ... ja
then ... take ... second t a little right is ... snout plane ... yes



Properties over Time

dann ... nehmen ... zweite t bisschen rechts ist ... rüssel hobel ... ja
then ... take ... second t a little right is ... snout plane ... yes

 yellow,T,TL,R1,C0



 blue,P,TL,R2,C

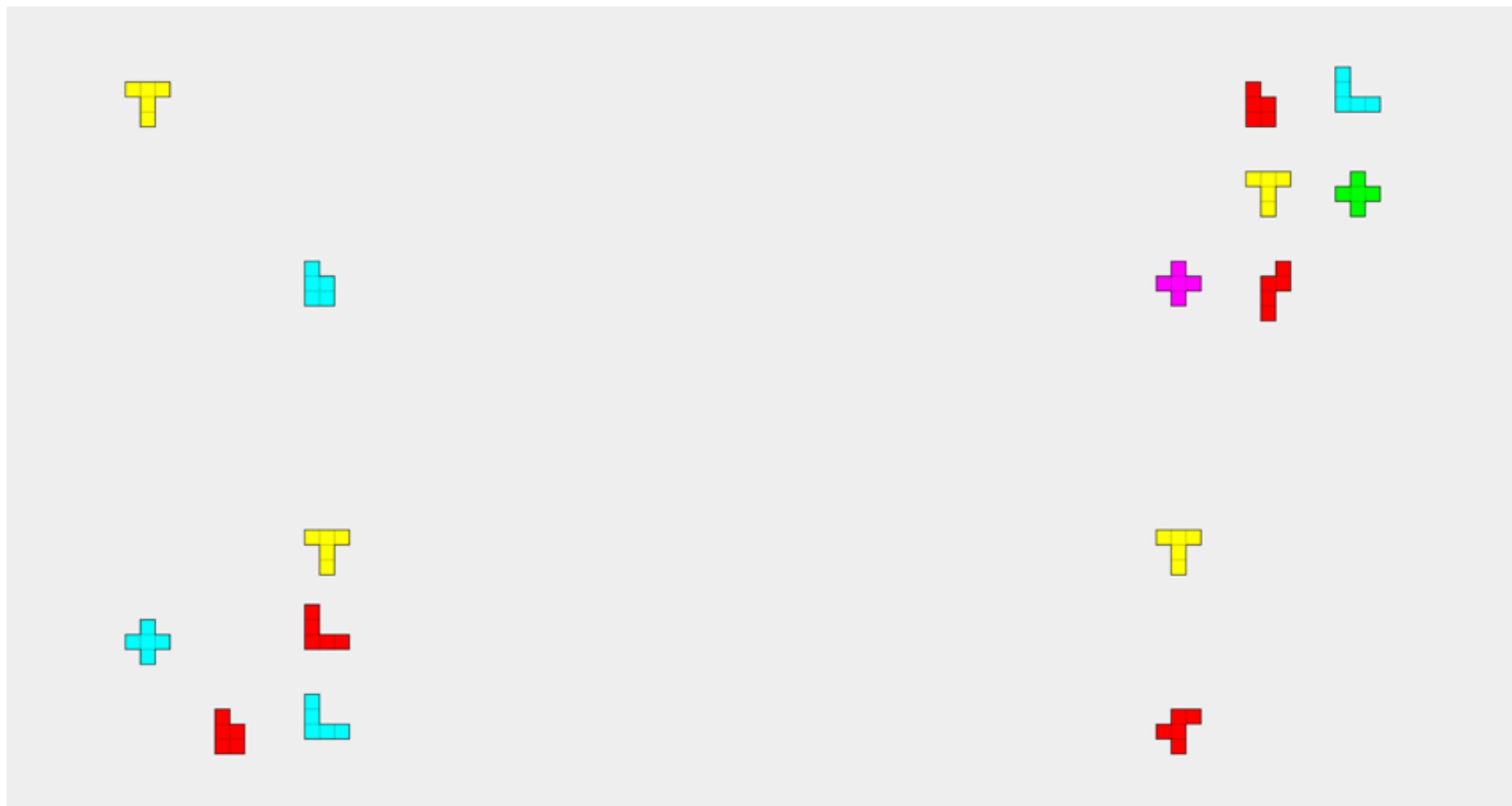


 yellow,T,R0,C2,**looked_at**

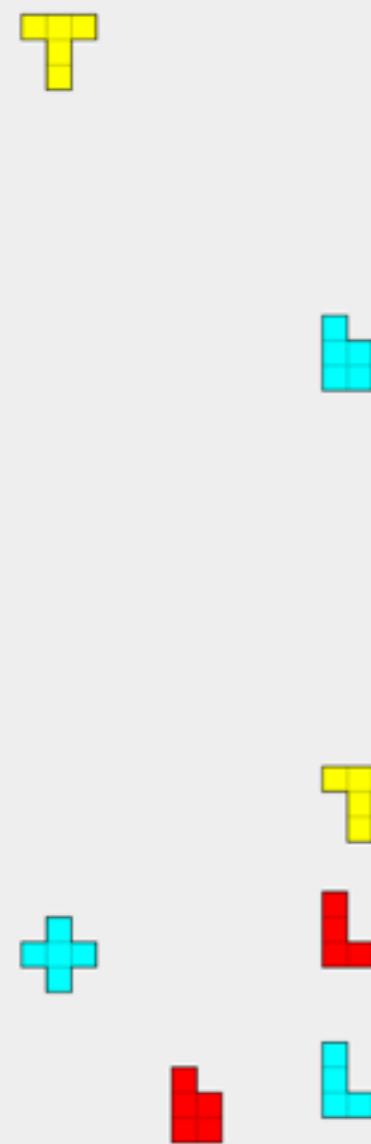


Properties over Time

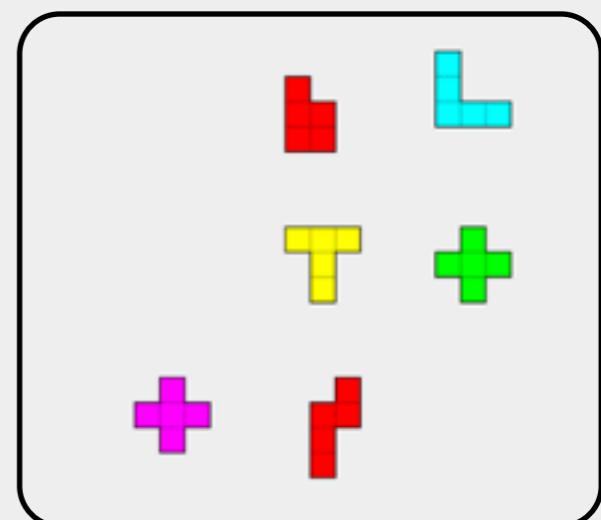
dann ... nehmen ... zweite t bisschen rechts ist ... rüssel hobel ... ja
then ... take ... second t a little right is ... snout plane ... yes



dann ... nehmen ... zweite t bisschen rechts ist ... rüssel hobel ... ja
then ... take ... second t a little right is ... snout plane ... yes



During pointing gesture,
all objects in grid receive
pointed_at property

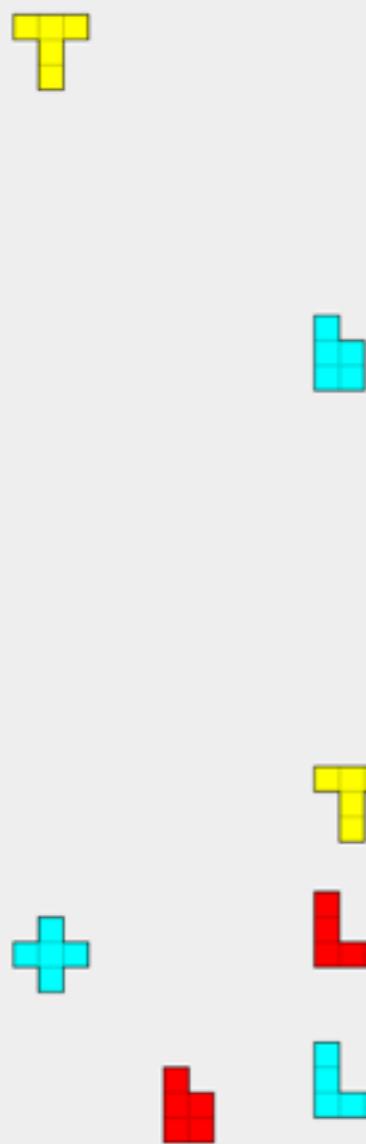


Properties over Time

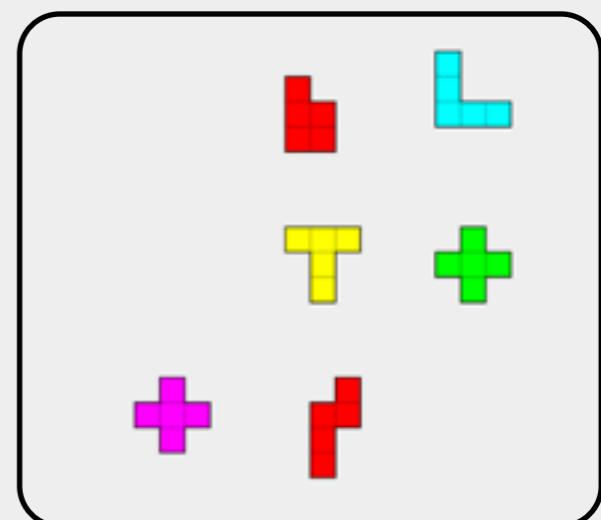
dann ... nehmen ... zweite t bisschen rechts ist ... rüssel hobel ... ja
then ... take ... second t a little right is ... snout plane ... yes



dann ... nehmen ... zweite t bisschen rechts ist ... rüssel hobel ... ja
then ... take ... second t a little right is ... snout plane ... yes

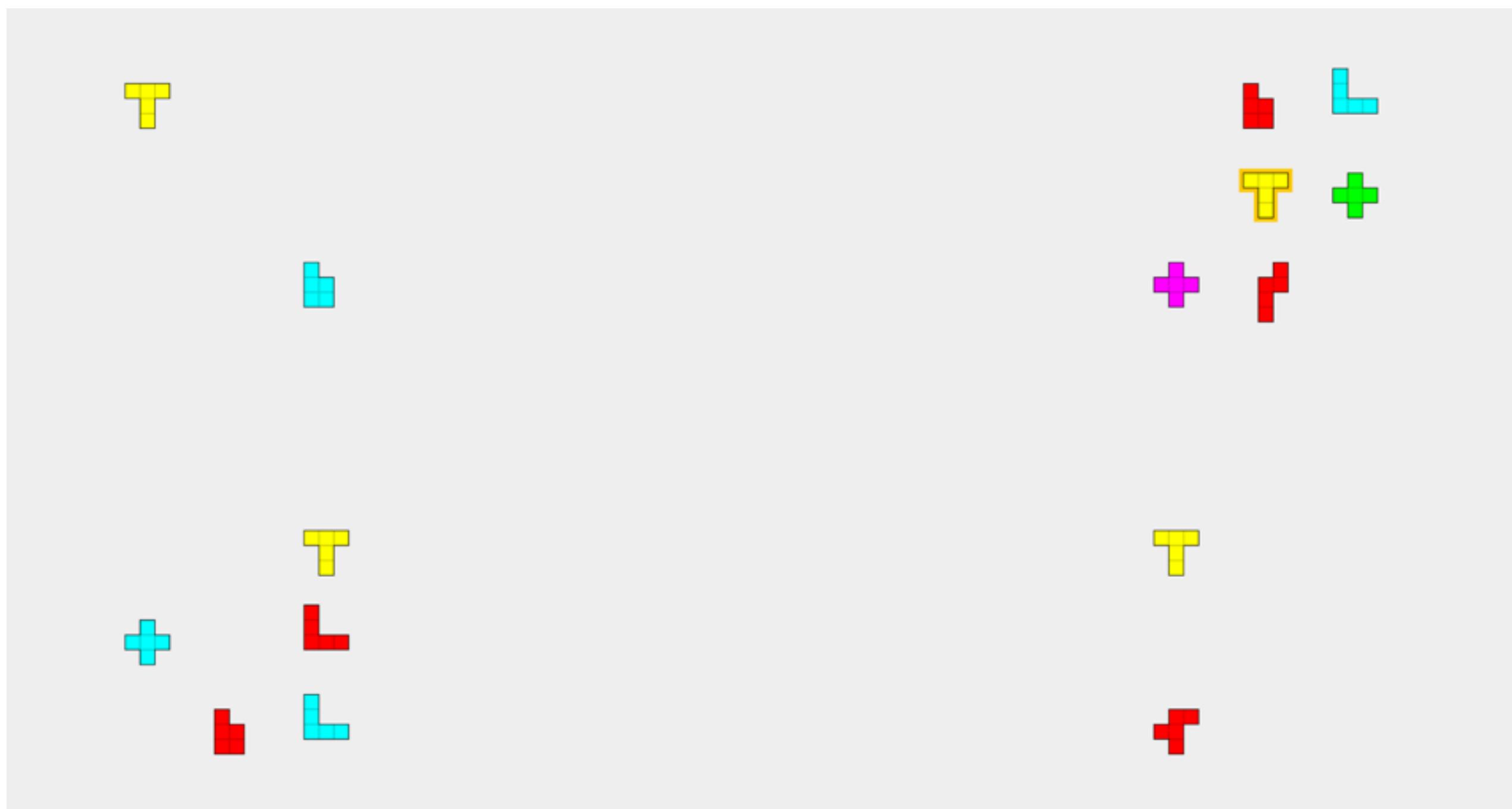


During pointing gesture,
all objects in grid receive
pointed_at property



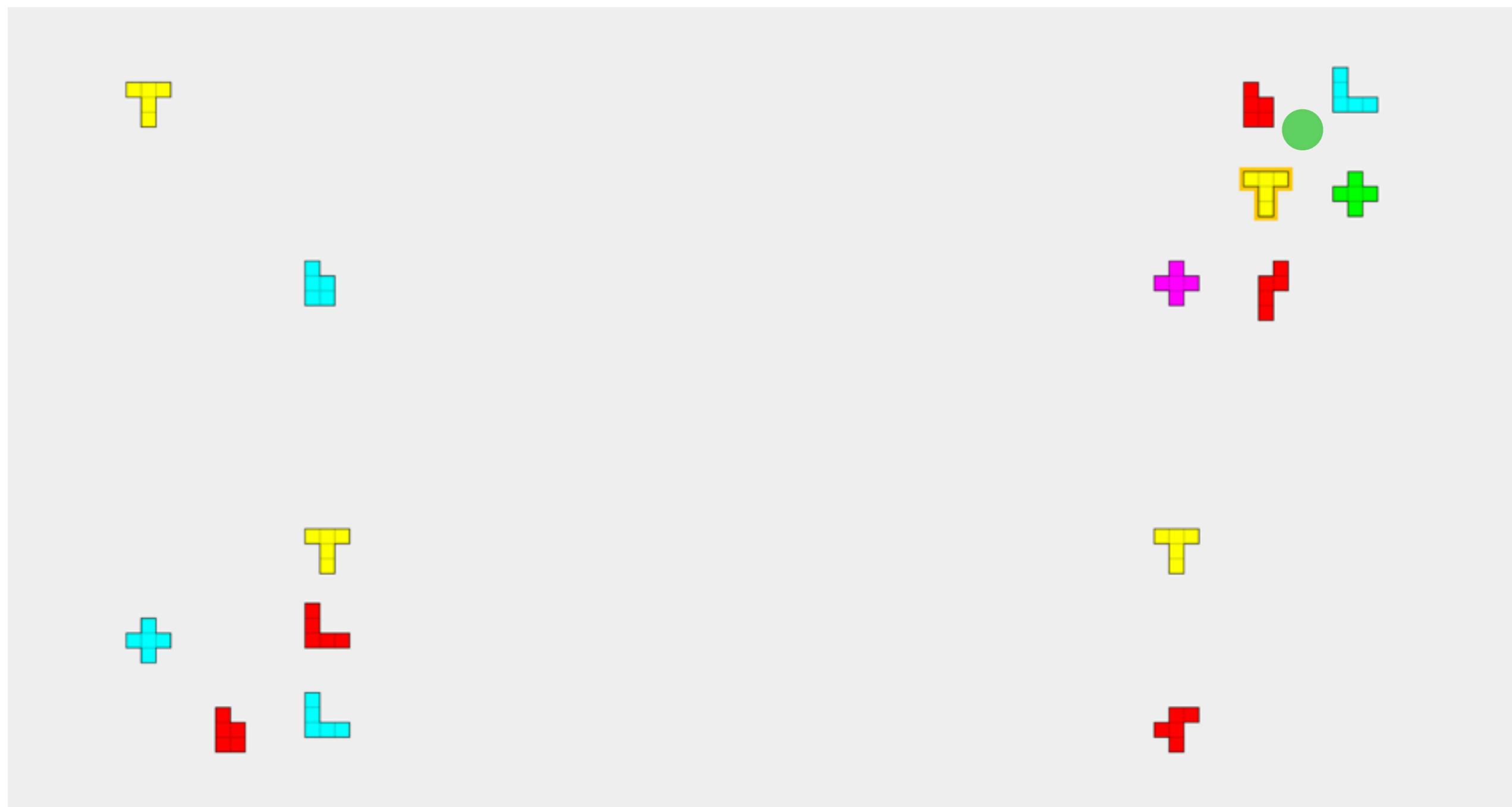
Properties over Time

dann ... nehmen ... zweite t bisschen rechts ist ... rüssel hobel ... ja
then ... take ... second t a little right is ... snout plane ... yes



Properties over Time

dann ... nehmen ... zweite t bisschen rechts ist ... rüssel hobel ... ja
then ... take ... second t a little right is ... snout plane ... yes



incremental statistical NLU

– Results for Pento-2013 Data –

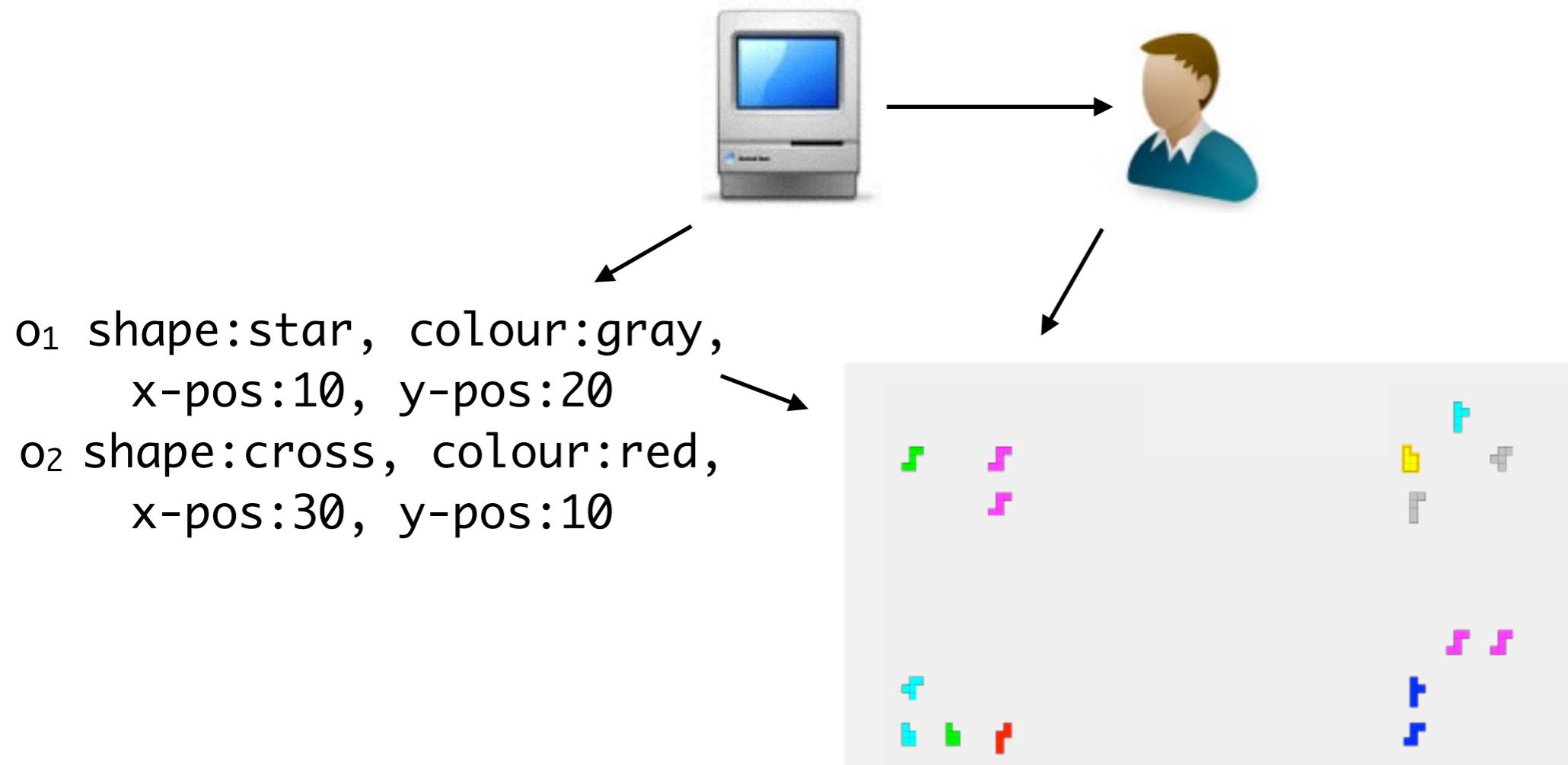


$$P(I|U, W) = \frac{1}{P(U|W)} P(I|W) \sum_R P(R|I, W) P(U|R)$$

Version	Acc	Top 2	Top 4
Gaze (baseline) NLU	18% 50%		
NLU + Gaze	53%	62%	80%
NLU + Point	52%	65%	90%
NLU + Gaze + Point	53%	70%	91%
NLU + Gaze-F	53%	65%	78%
NLU + Point-F	57%	68%	88%
NLU+Gaze-F+Point-F	56%	69%	85%

incremental statistical NLU

– The Pento-2013 Data –



incremental statistical NLU

– The Pento-2013 Data, via CV –

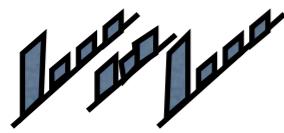
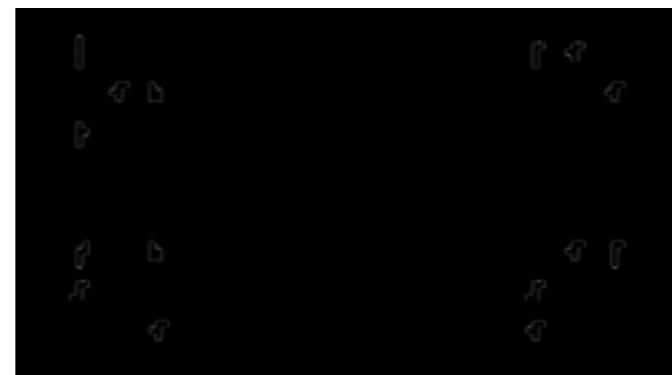
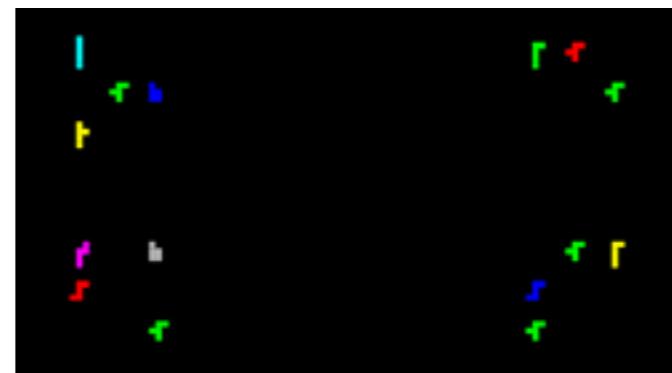


joint work with Liva Dia

incremental statistical NLU

– The Pento-2013 Data, via CV –

o₁ shape:star, colour:gray,
x-pos:10, y-pos:20
o₂ shape:cross, colour:red,
x-pos:30, y-pos:10,
. . .

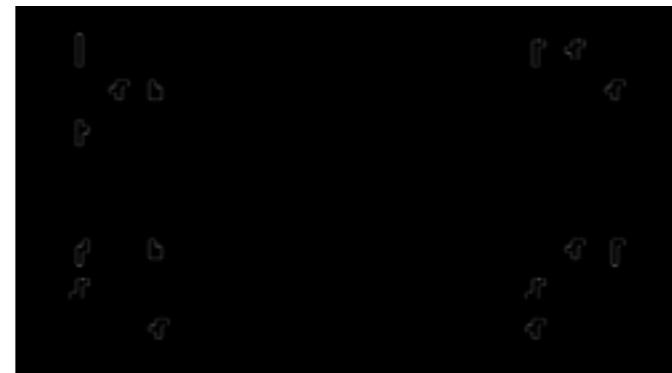
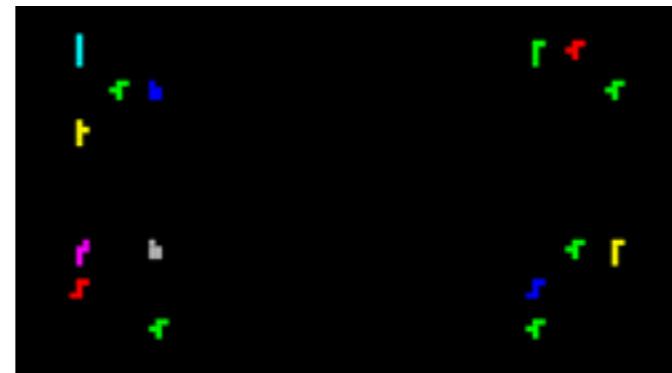


...

incremental statistical NLU

– The Pento-2013 Data, via CV –

o₁ shape:star, colour:gray,
x-pos:10, y-pos:20
o₂ shape:cross, colour:red,
x-pos:30, y-pos:10,
. . .



o₁ shape:star, colour:gray,
x-pos:10, y-pos:20
o₂ shape:cross, colour:red,
x-pos:30, y-pos:10,
. . .

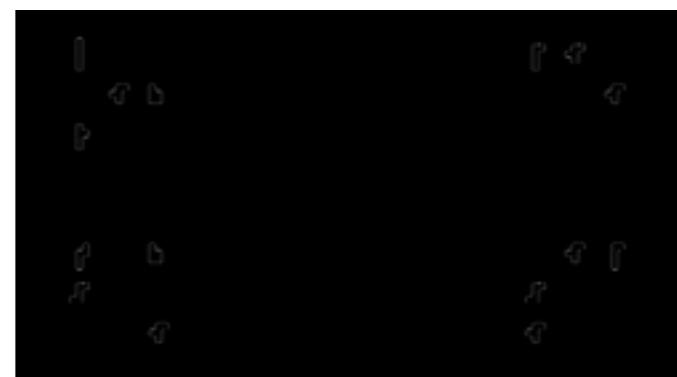
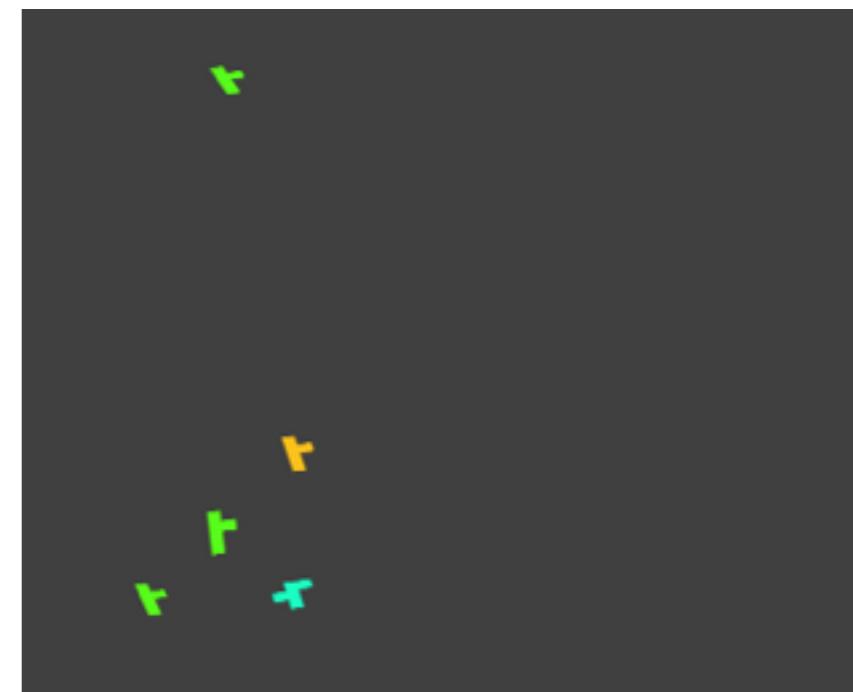


. . .

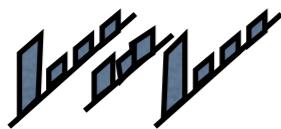
incremental statistical NLU

– The Pento-2013 Data, via CV –

```
o1 shape:star, colour:gray,  
    x-pos:10, y-pos:20  
o2 shape:cross, colour:red,  
    x-pos:30, y-pos:10,  
    . . .
```



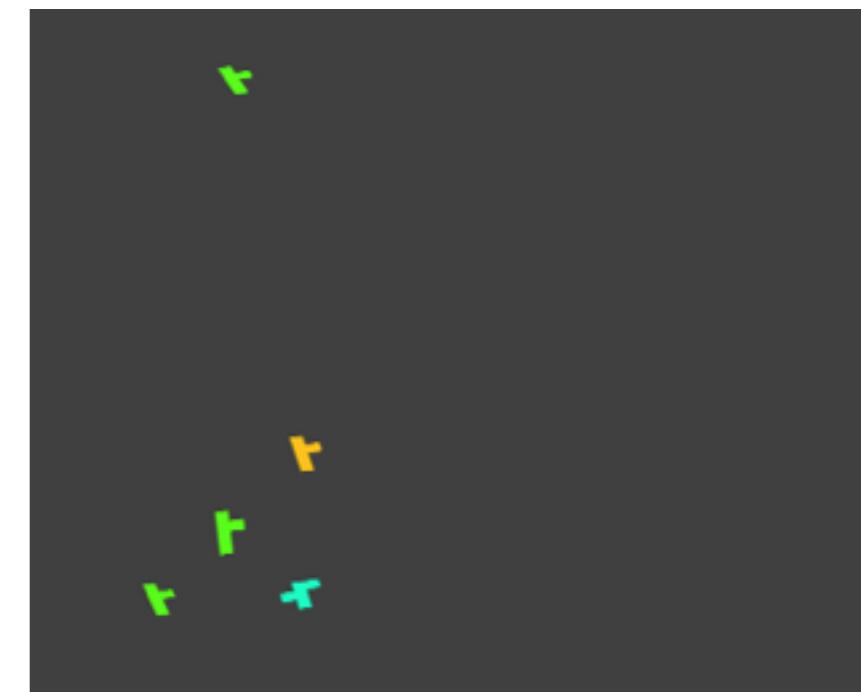
...



incremental statistical NLU

– The Pento-2013 Data, via CV –

	o	o
star	1	0
cross	0	1
gray	1	0
red	0	1



	o	o
star	0.8	0.1
cross	0.2	0.9
gray	0.9	0.4
red	0.1	0.6



	o	o
star	0.4	0.05
cross	1	0.45
gray	0.45	0.2
red	0.05	0.3



...



incremental statistical NLU

– The Pento-2013 Data, via CV –

	o	o
star	1	0
cross	0	1
gray	1	0
red	0	1



	o	o
star	0.8	0.1
cross	0.2	0.9
gray	0.9	0.4
red	0.1	0.6



	o	o
star	0.4	0.05
cross	1	0.45
gray	0.45	0.2
red	0.05	0.3



Version	Acc
Base	0.78
CV orig, argmax	0.77
CV orig, distr	0.75

...

incremental statistical NLU

– The Pento-2013 Data, via CV –

	o	o
star	1	0
cross	0	1
gray	1	0
red	0	1



	o	o
star	0.8	0.1
cross	0.2	0.9
gray	0.9	0.4
red	0.1	0.6



	o	o
star	0.4	0.05
cross	1	0.45
gray	0.45	0.2
red	0.05	0.3



Version	Acc
Base	0.78
CV orig, argmax	0.77
CV orig, distr	0.75
CV filt, argmax	0.59
CV filt, distr	0.60

...

incremental statistical NLU

– generative vs. discriminative –

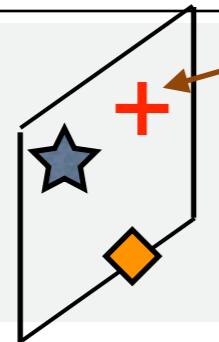
res.

x=02

P(I | S, W)

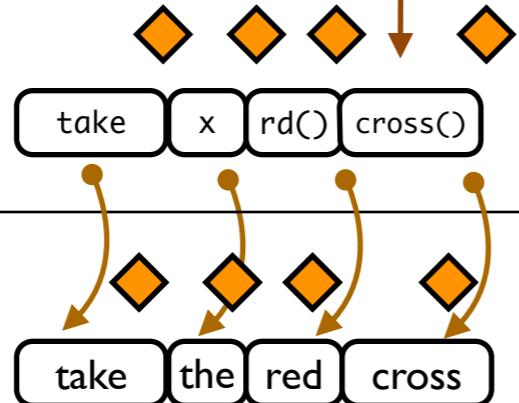
$$P(I | S, W) = \frac{P(S | I, W) P(I|W)}{P(S|W)}$$

vis

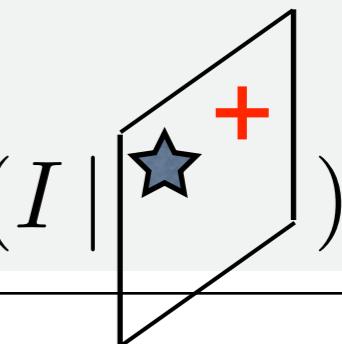


$$P(I | S, W) = P_W(I | \star)$$

sem



P(I | S, W)



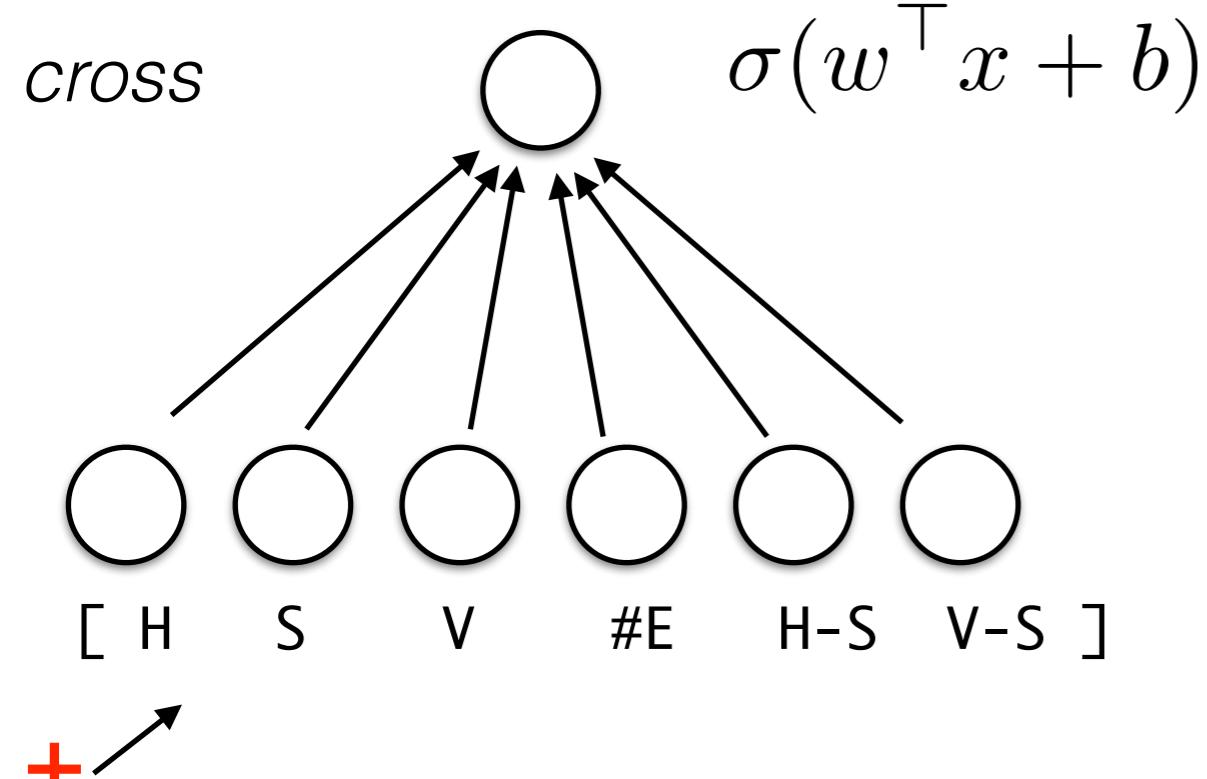
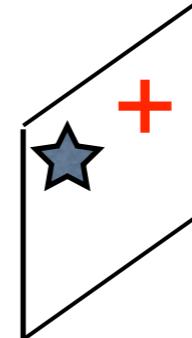
word
rec.



incremental statistical NLU

– discriminative model –

$$P(I | S, W) = P_W(I | \boxed{\star})$$



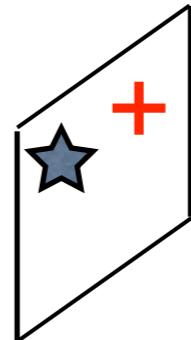
Training:

- * for each occurrence of word in corpus:
 - * pair it with (features of) object utterance referred to [positive example]
 - * pair with (features of) n objects in scene it didn't refer to [negative examples]
- * train (binary) logistic regression classifier for word

incremental statistical NLU

– discriminative model –

$$P(I | S, W) = P_W(I | \boxed{\star})$$



Testing:

- * for each word of utterance:
 - * test each object in scene against classifier of this word
 - * update distribution from previous step with distribution from this step

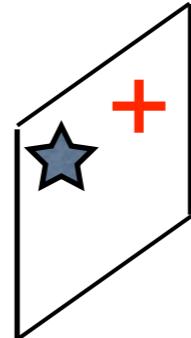
Training:

- * for each occurrence of word in corpus:
 - * pair it with (features of) object utterance referred to [positive example]
 - * pair with (features of) n objects in scene it didn't refer to [negative examples]
- * train (binary) logistic regression classifier for word

incremental statistical NLU

– discriminative model –

$$P(I | S, W) = P_W(I | \boxed{\star})$$



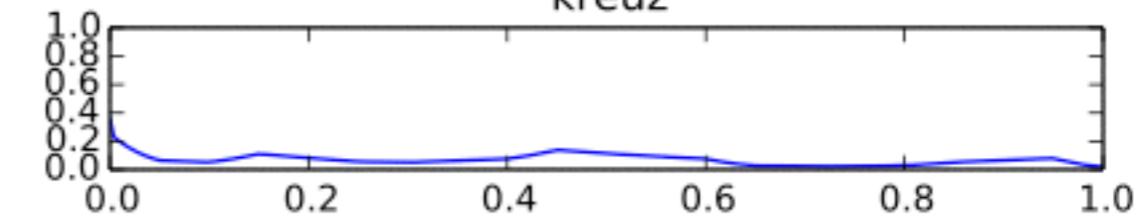
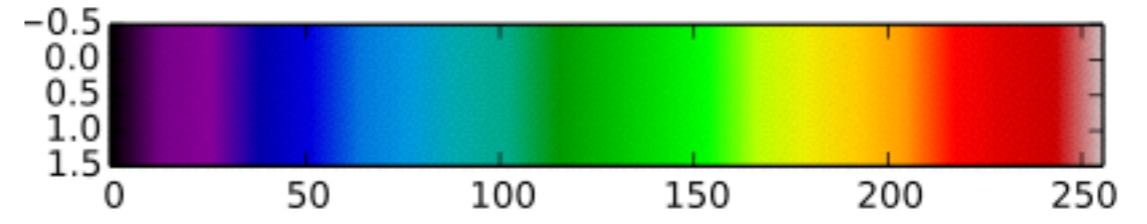
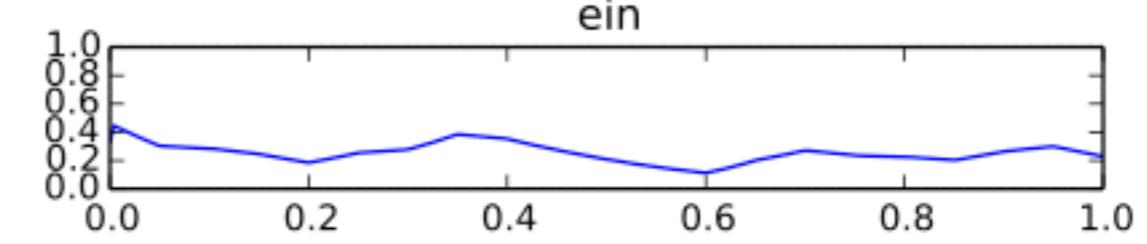
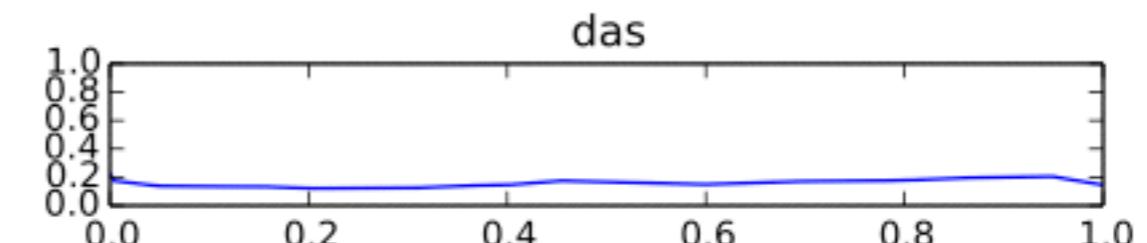
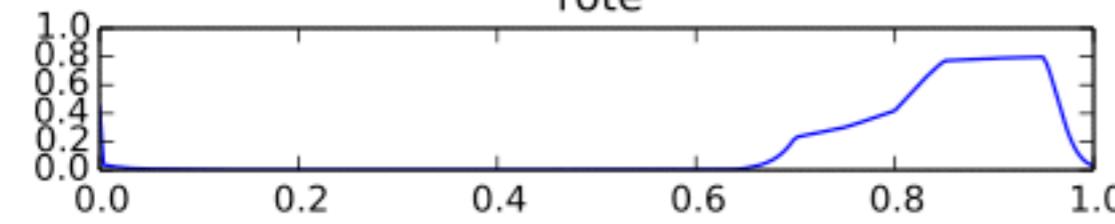
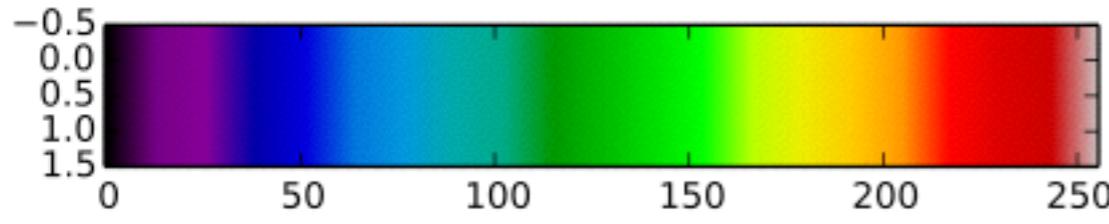
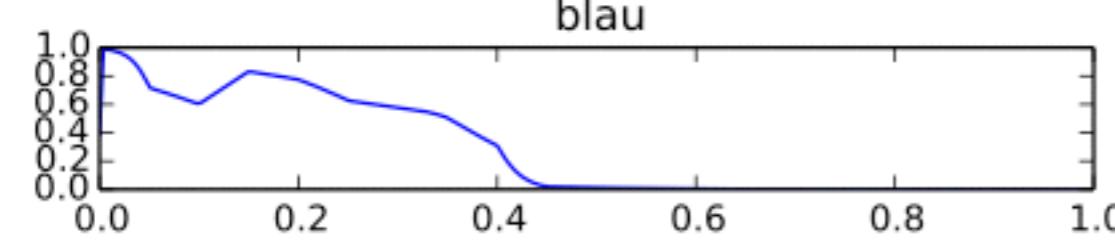
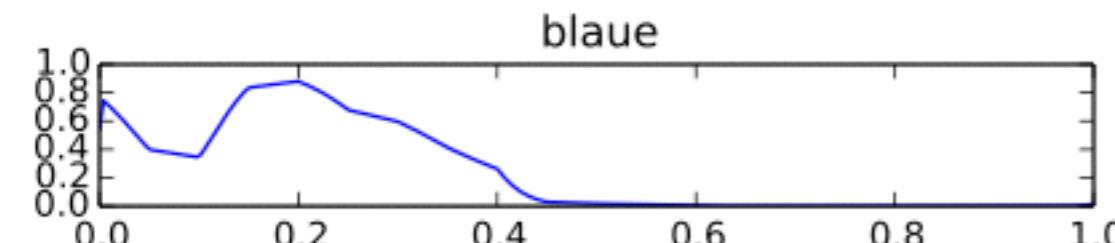
Training:

- * for each occurrence of word in corpus:
 - * pair it with (features of) object utterance referred to [positive example]
 - * pair with (features of) n objects in scene it didn't refer to [negative examples]
 - * train (binary) logistic regression classifier for word

Version	Acc
Base	0.78
CV orig, argmax	0.77
CV orig, distr	0.75
CV filt, argmax	0.59
CV filt, distr	0.60
discr, CV filt & distr.	0.66

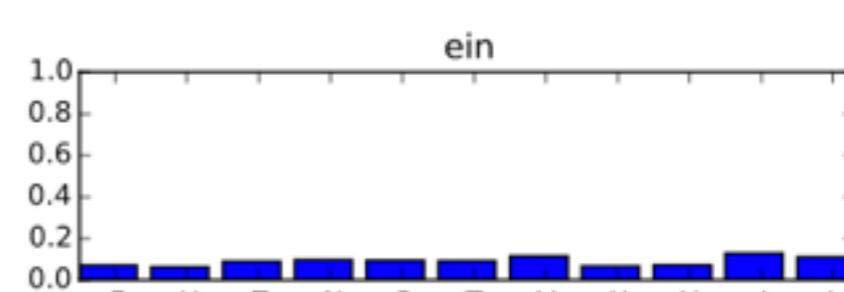
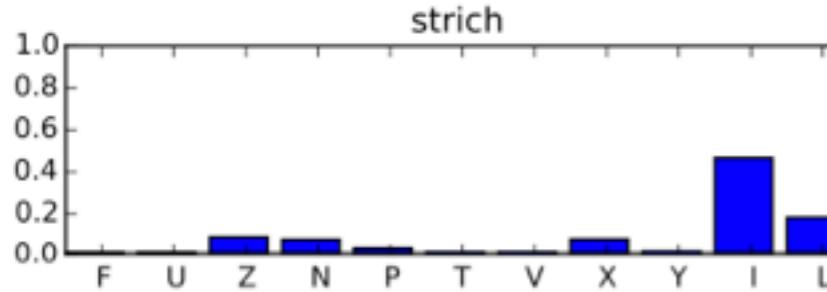
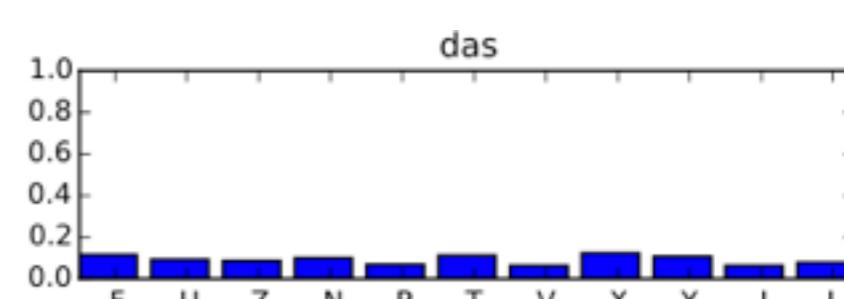
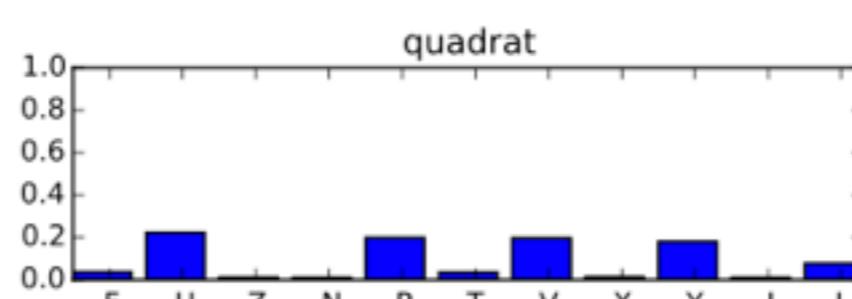
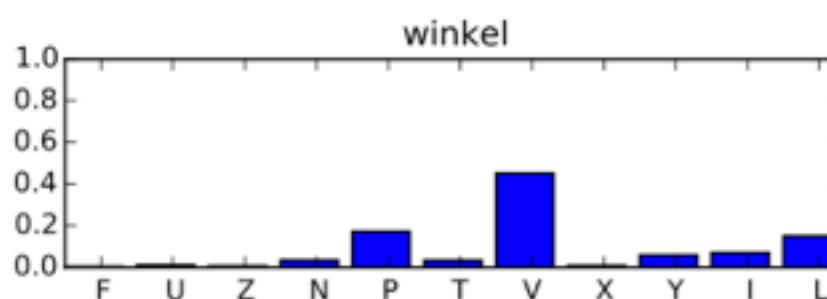
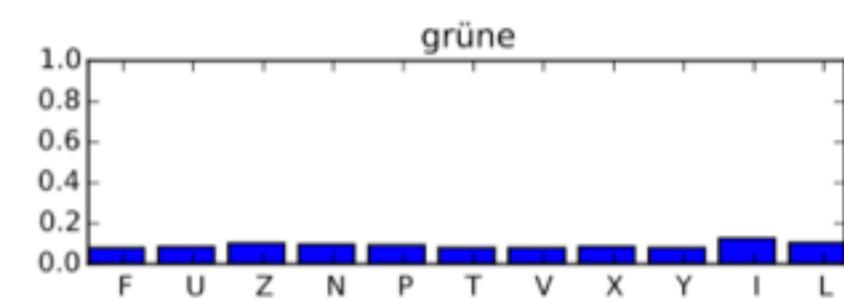
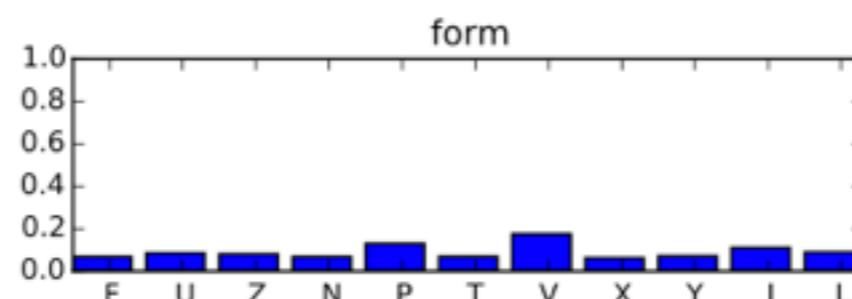
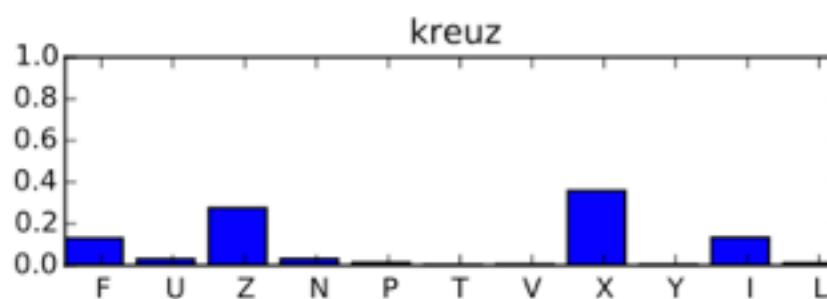
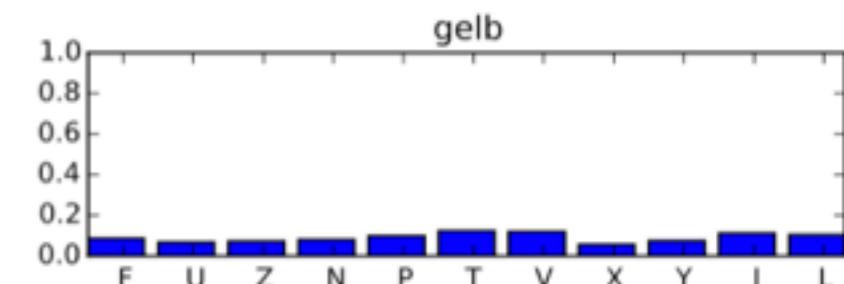
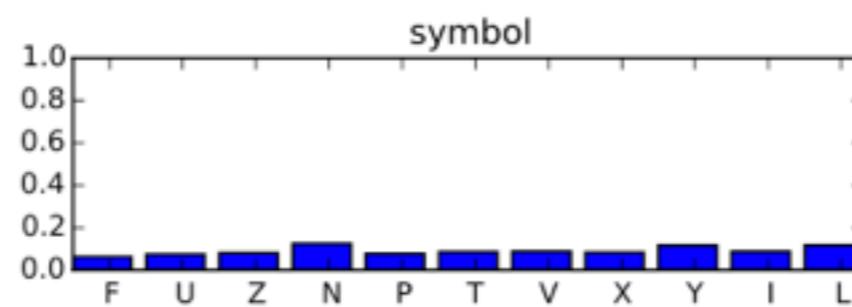
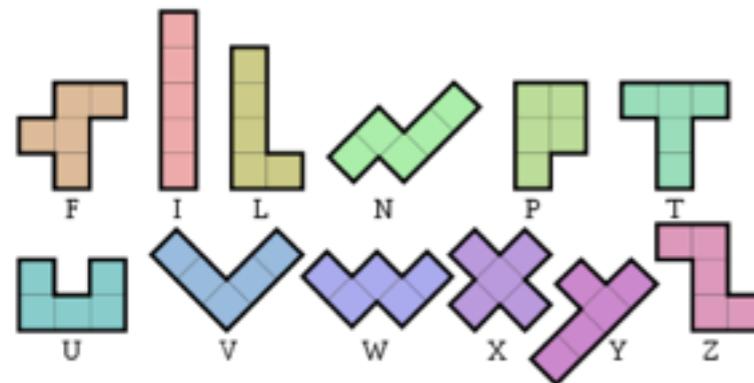
incremental statistical NLU

– discriminative model.. word semantics? –



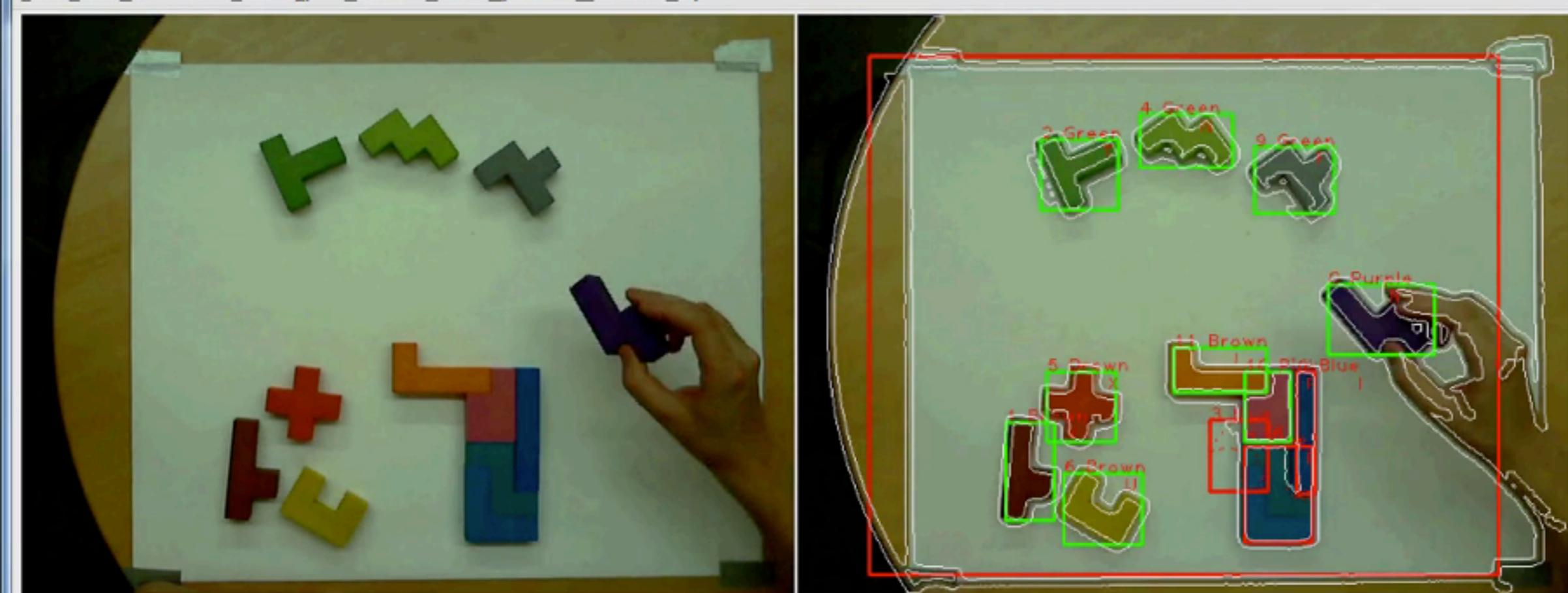
incremental statistical NLU

– discriminative model.. word semantics? –



incremental statistical NLU

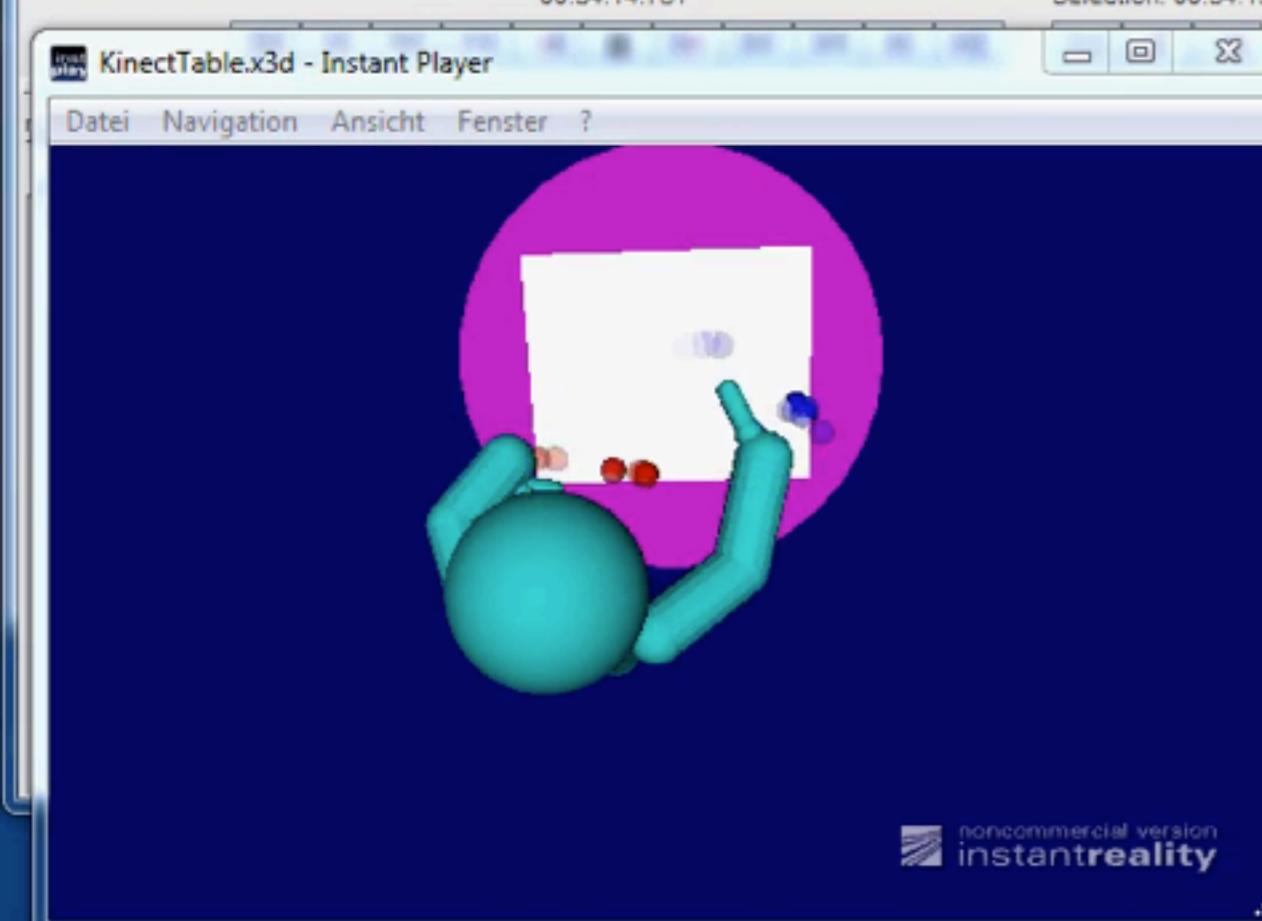
- models that learn from annotated data (instances of reference), and that
- can be applied incrementally, on-line, and continuously produce hypotheses,
- can make use of additional information sources such as gestures, eye gaze
- can work with non-symbolic input (images)
- perform relatively well



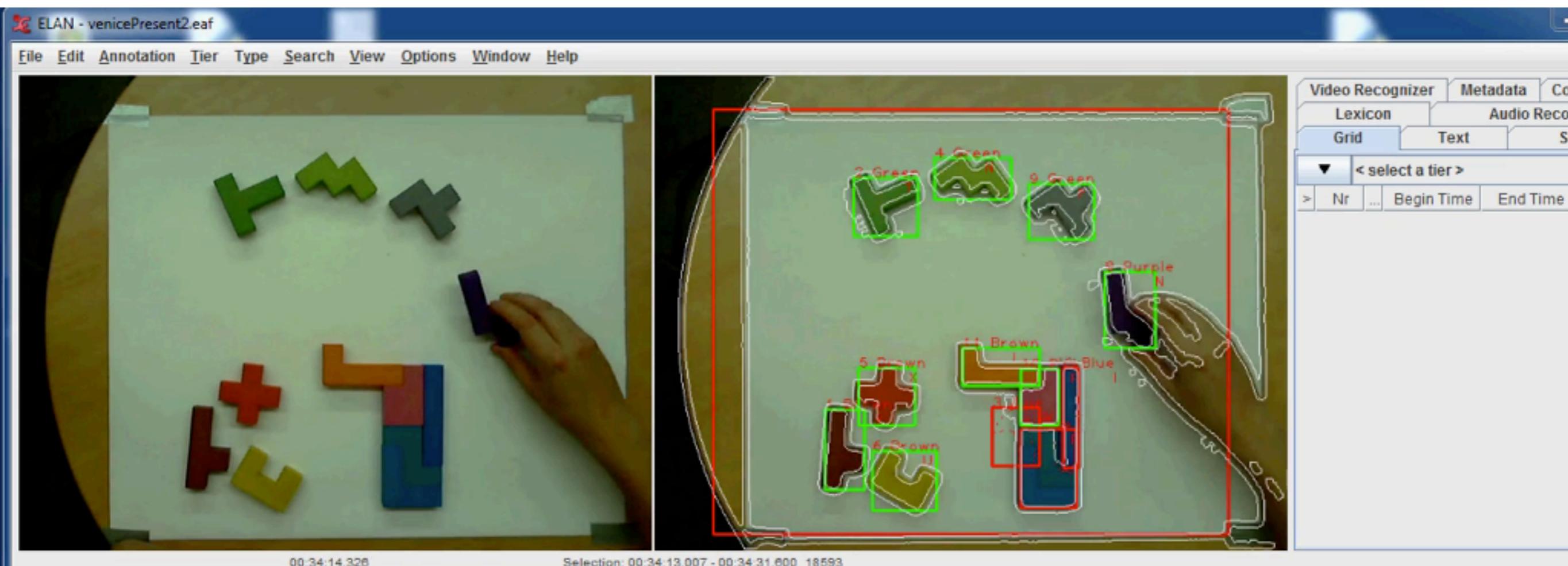
Video Recognizer
Lexicon
Grid
Text

< select a tier >

Nr ... Begin Time



```
C:\Windows\system32\cmd.exe - java -jar VeniceHub.jar @..\config_debug.args
Progress of last skip: 1,00 overall < lokal 1,00>
skip no. 3: 7731 bytes of data <avgBytePerS=23572>
skip no. 0: 396012 bytes of data <avgBytePerS=14423>
skip no. 0: 396012 bytes of data <avgBytePerS=14423>
skip no. 0: 404637 bytes of data <avgBytePerS=14423>
skip no. 0: 419031 bytes of data <avgBytePerS=14423>
skip no. 0: 433468 bytes of data <avgBytePerS=14423>
skip no. 0: 447891 bytes of data <avgBytePerS=14423>
skip no. 0: 462329 bytes of data <avgBytePerS=14423>
skip no. 0: 476781 bytes of data <avgBytePerS=14423>
skip no. 0: 491204 bytes of data <avgBytePerS=14423>
skip no. 0: 505612 bytes of data <avgBytePerS=14423>
skip no. 0: 520064 bytes of data <avgBytePerS=14423>
skip no. 0: 534473 bytes of data <avgBytePerS=14423>
skip no. 0: 549011 bytes of data <avgBytePerS=14423>
skip no. 0: 563448 bytes of data <avgBytePerS=14423>
skip no. 0: 577857 bytes of data <avgBytePerS=14423>
skip no. 0: 592367 bytes of data <avgBytePerS=14423>
skip no. 0: 606804 bytes of data <avgBytePerS=14423>
skip no. 0: 621227 bytes of data <avgBytePerS=14423>
skip no. 0: 635650 bytes of data <avgBytePerS=14423>
skip no. 0: 650073 bytes of data <avgBytePerS=14423>
skip no. 0: 664496 bytes of data <avgBytePerS=14423>
reset
resetting
```



00:34:14.326

Selection: 00:34:13.007 - 00:34:31.600 18593

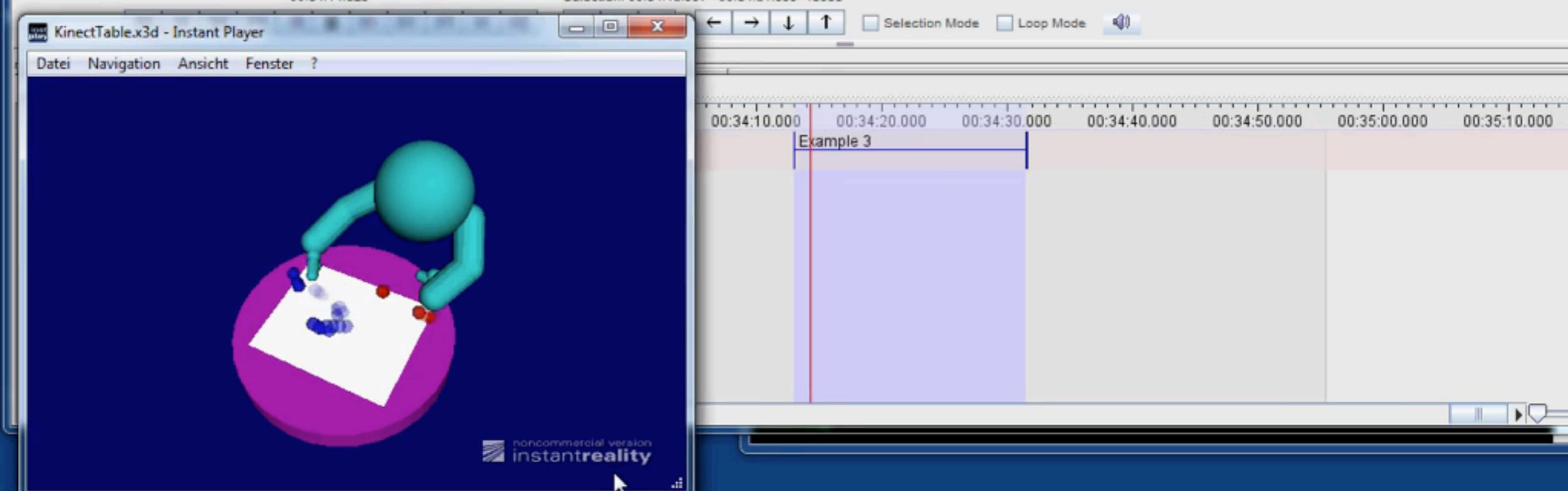
Video Recognizer Metadata Co

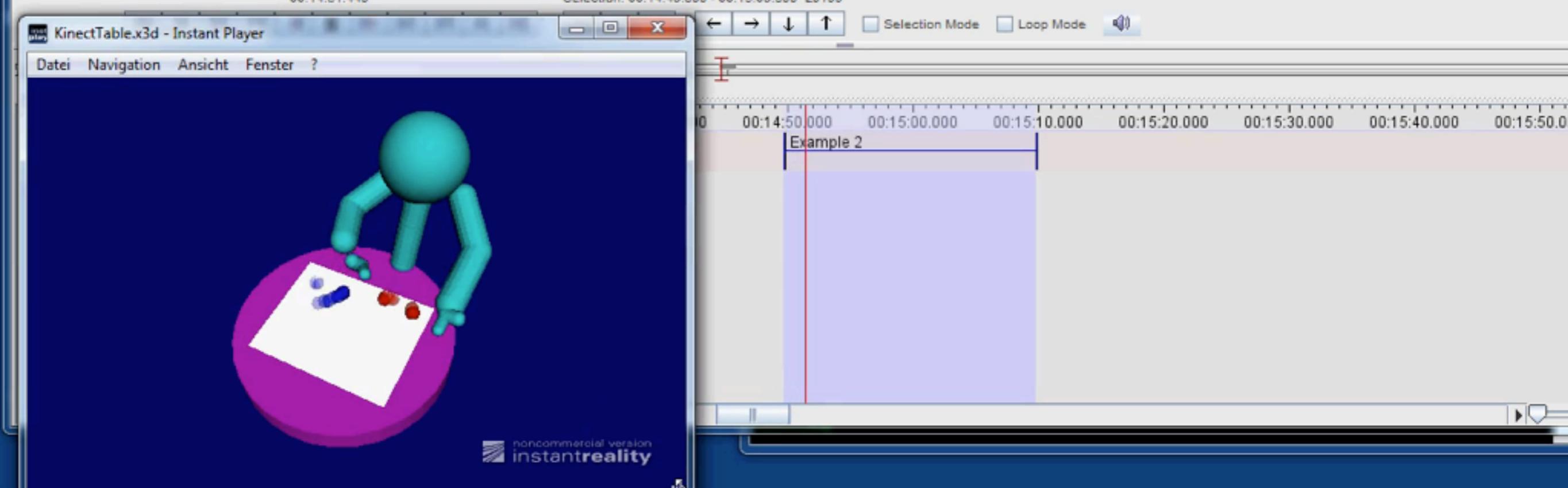
Lexicon Audio Reco

Grid Text S

< select a tier >

Nr ... Begin Time End Time





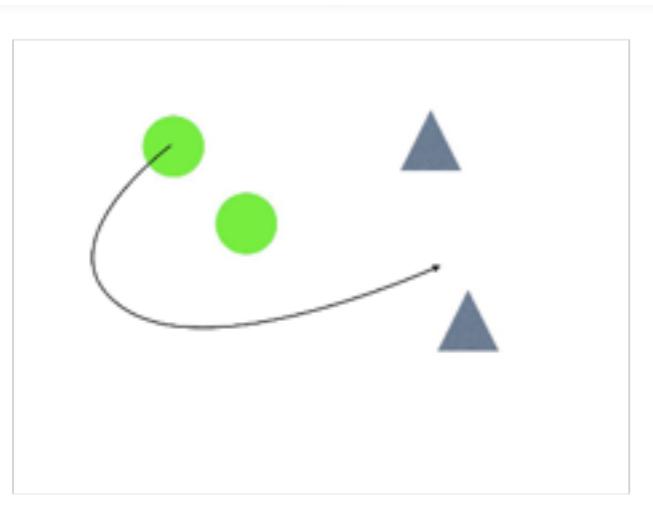
incremental statistical NLU

- models that learn from annotated data (instances of reference), and that
- can be applied incrementally, on-line, and continuously produce hypotheses,
- can make use of additional information sources such as gestures, eye gaze
- can work with non-symbolic input (images)
- perform relatively well

towards a model of sit. com.

– future work –

- learning through interaction / active supervision
- compositionality
- more types of gestures
- managing uncertainty; clarification sequences; active grounding
- embodiment



hier ist ein graues Dreieck
here is a gray triangle



und hier ist ein grüner Kreis
and here is a green circle



hier ist noch ein grüner Kreis
here is another green circle



und hier ist noch ein graues Dreieck
and here is another gray triangle



und von dem oberen grünen Kreis [...]
and from the top green cross [...]

towards a model of sit. com.

– future work –

- learning through (active) interaction / supervision
- compositionality
- more types of gestures
- managing uncertainty; clarification sequences; active grounding
- embodiment

Overview

- **Part I: Foundations**

- coordination, convention
- communicative intentions
- non-conventional meaning
- grounding
- turn-taking
- disfluencies

- **Part II: Computational Models**

- approaches to dialogue modelling
- incremental processing, turn-taking
- an example: grounded semantics

2006-2012

Incremental
Processing and
Projection in
Dialogue (InPro)



2010-2014

SFB (collaborative
research center)
673, alignment in
communication



University of Potsdam

2010-2018

Center of
Excellence,
Cognitive
Interaction
Technology (CIT-
EC)



Bielefeld University

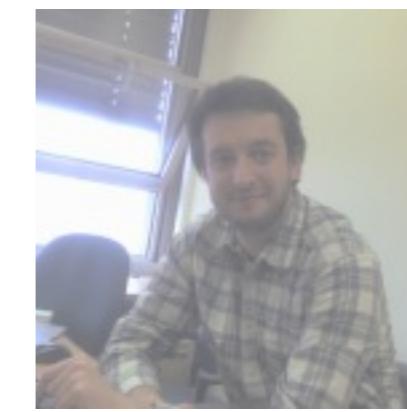
2014-2016

Disfluencies,
Exclamations, &
Laughter in
Dialogue



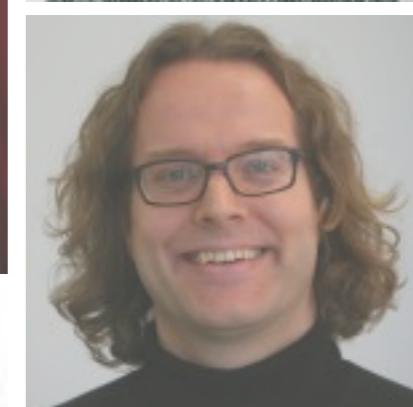
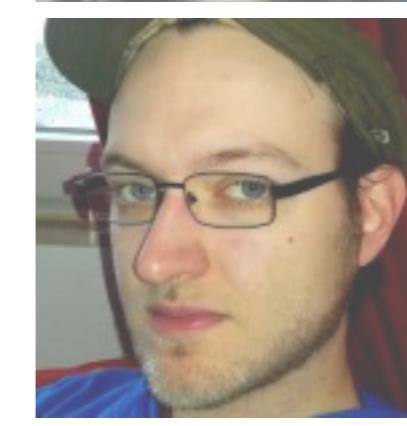
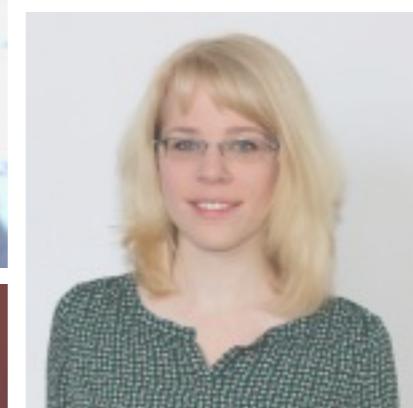
- Post-Docs

- Spyros Kousidis (PhD Dublin)
- Iwan de Kok (PhD U Twente)
- Julian Hough (PhD Queen Mary, U London)



- PhD Students

- Casey Kennington
- Ting Han
- Birte Carlmeyer
- Simon Betz



- Alumni

- Timo Baumann (PhD)
- Gabriel Skantze (Post-Doc)
- Okko Buß (PhD)
- Michaela Atterer (Post-Doc)



thank you!