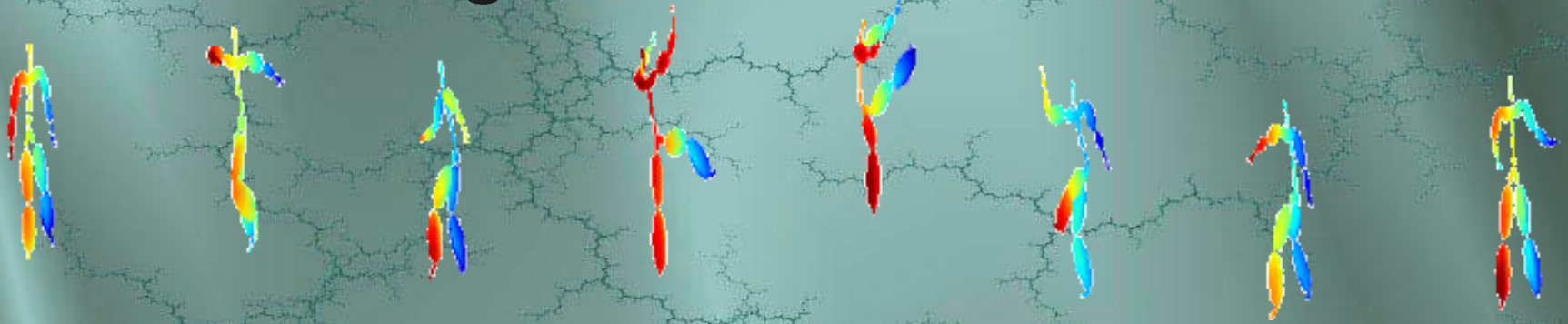


Human Activity Language: Grounding Concepts with a Linguistic Framework



Gutemberg Guerra-Filho

{guerra@cs.umd.edu}

Yiannis Aloimonos

{yiannis@cfar.umd.edu}

Computer Vision Laboratory
Department of Computer Science
University of Maryland



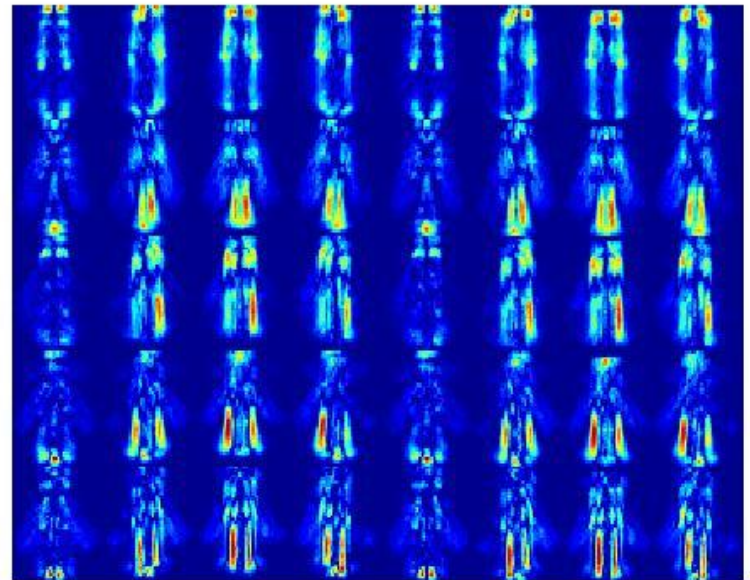
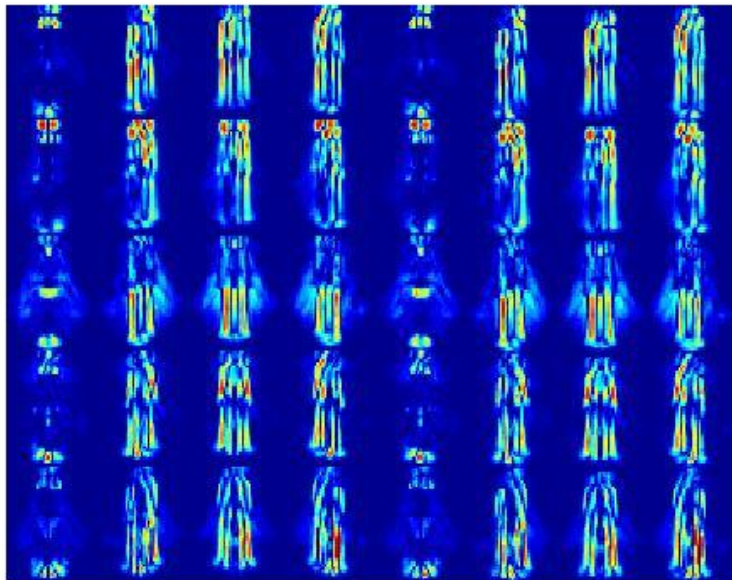
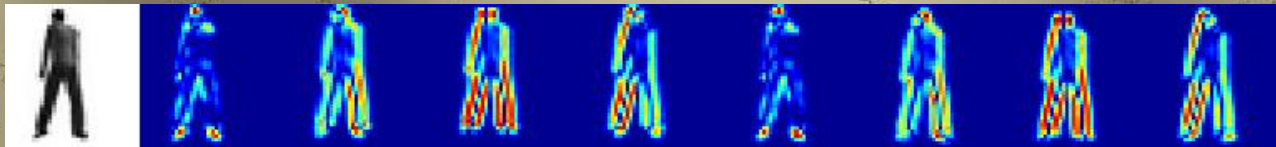
Theme

- To close the semantic gap in multimedia technologies, we need to **understand human action**
- There are at least 3 spaces devoted to human action: The Visual, the Motoric, and the Language Space.
- Each of these spaces is characterized by a distinct language, with its own alphabet, words, and syntax.

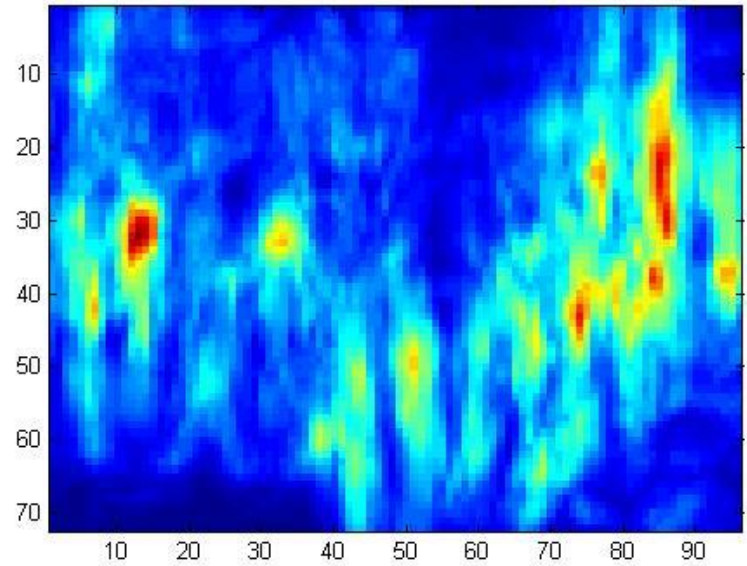
Initial Meeting

- Multimedia
- Semantics
- Semantics arises from human action
- Big brother problem

VHF's: Visual Human Filters

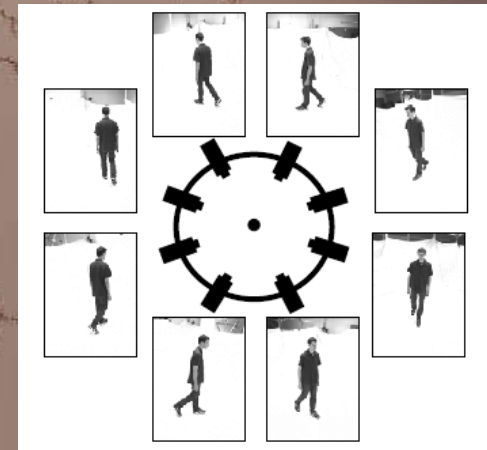
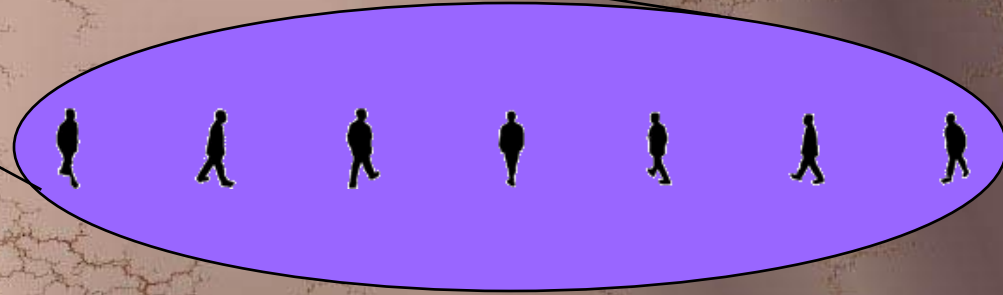


Applying the VHF's



Visual Approach: Sequences of poses

p_1 Stand	
p_2 Bent Knees	
p_3 Legs Apart(1)	
p_4 Legs Together	
p_5 Legs Apart(2)	
p_6 Kick Leg Behind	
p_7 Kick Leg Front	
p_8 Kick Legs Together	
p_9 Kneel	
p_{10} Half Squat Down	
p_{11} Squat	
p_{12} Half Squat Up	
p_{13} Half Bend Down	
p_{14} Full Bend	
p_{15} Half Bend Up	
p_{16} Start Sit Down	
p_{17} Half Sit Back	
p_{18} Full Sit	
p_{19} Half Sit Front	
p_{20} Start Sit Up	



What are “Key” poses?

- Extremal poses of the body.
- How are they found?
- Single-view example:



Video

Horizontal motion

Red = right

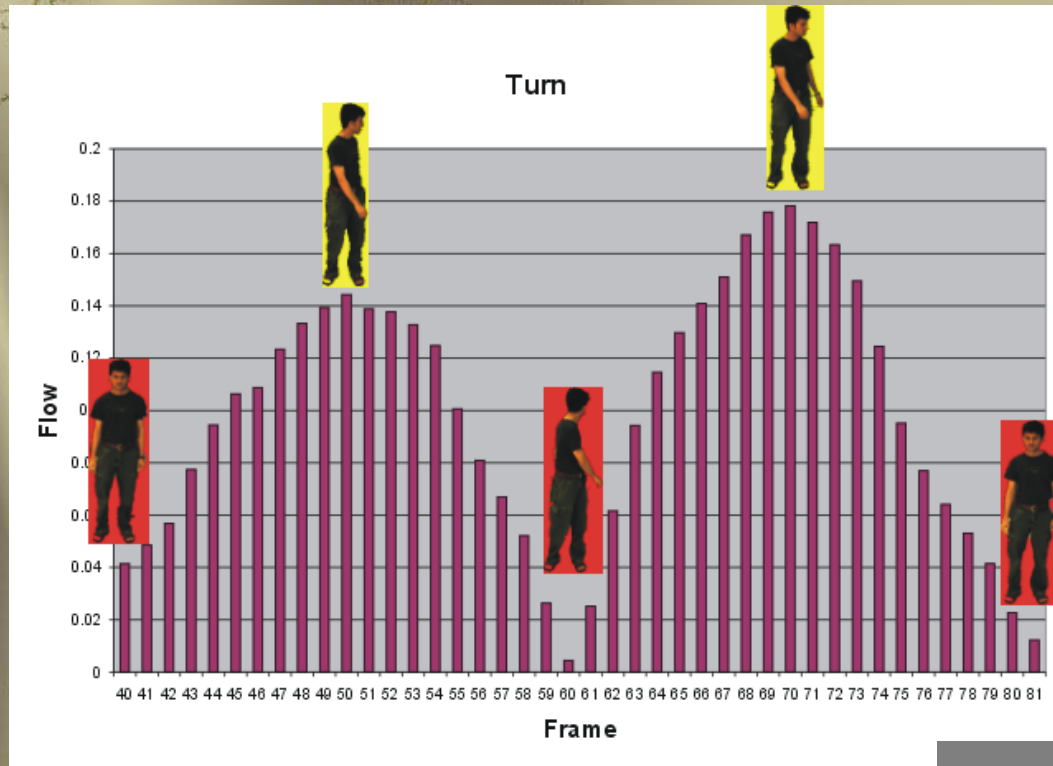
Blue = left

Vertical motion

Red = Down

Blue = Up

What are “Key” poses?



← Mean
flow magnitude
(in person's
reference frame)

Key frames →



Pose Grammar

- Probabilistic context-free Grammar (PCFG).

$$Start \rightarrow V \quad p = 1$$

$$V \rightarrow VA \mid A \quad p = \frac{1}{2}$$

$$A \rightarrow A_1 \mid A_2 \mid \dots \mid A_g \quad \forall i, p(A_i \mid A) = 1/g$$

$$A_i \rightarrow q_{ab} q_{bc} q_{cd} \dots \quad p(q_{ab} q_{bc} q_{cd} \dots \mid A_i) = 1$$

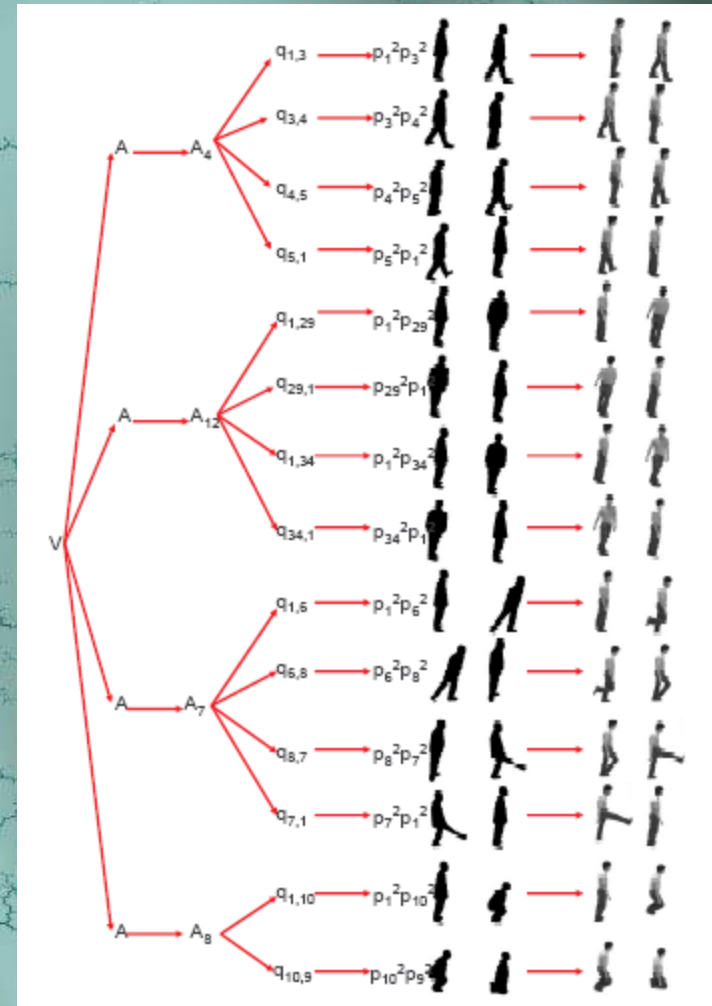
$$q_{cd} \rightarrow p_c^u p_d^v \quad \sum_{\substack{\text{allowed} \\ u,v}} p(p_c^u p_d^v \mid q_{cd}) = 1$$

$$p_i^v \rightarrow s_k \quad p(s_k \mid p_i^v) \text{ obtained at runtime}$$

Rules created
from training
data

Parse an input video

1. Key frame detection.
2. Silhouette matching on keyframes.
3. Computation of $P(s_k | p_i^v)$ as shown earlier.
4. Probabilistic parsing using the PCFG.



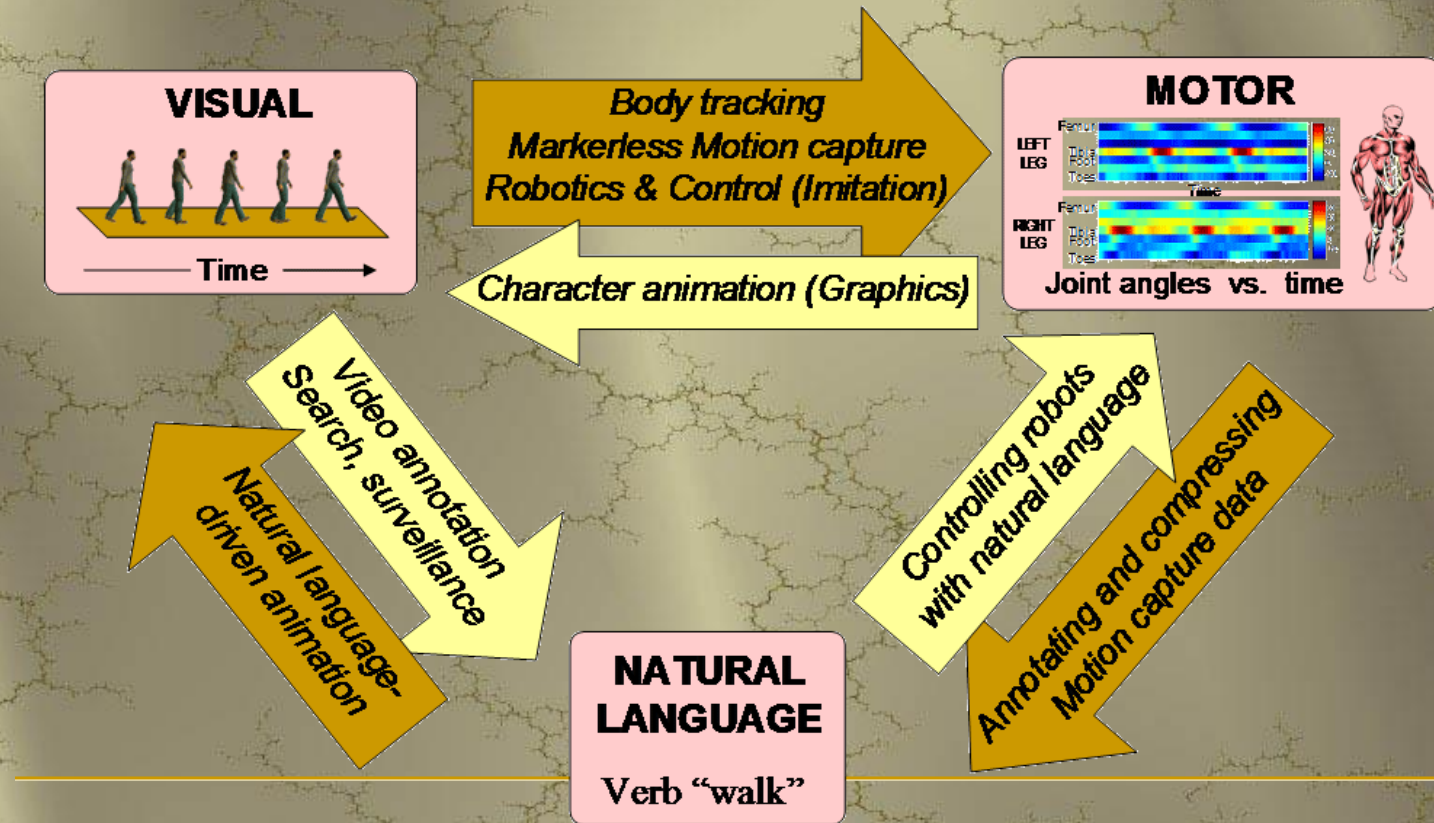
The background of the image is a complex, golden-brown texture. It features a network of fine, irregular cracks that create a marbled or cracked-glass effect. The overall color palette is warm, ranging from light beige to deep, dark brown tones, with a prominent golden sheen. The lighting appears to come from the left, creating a subtle gradient and highlighting the intricate patterns of the texture.

One year later...

- That's What I Found

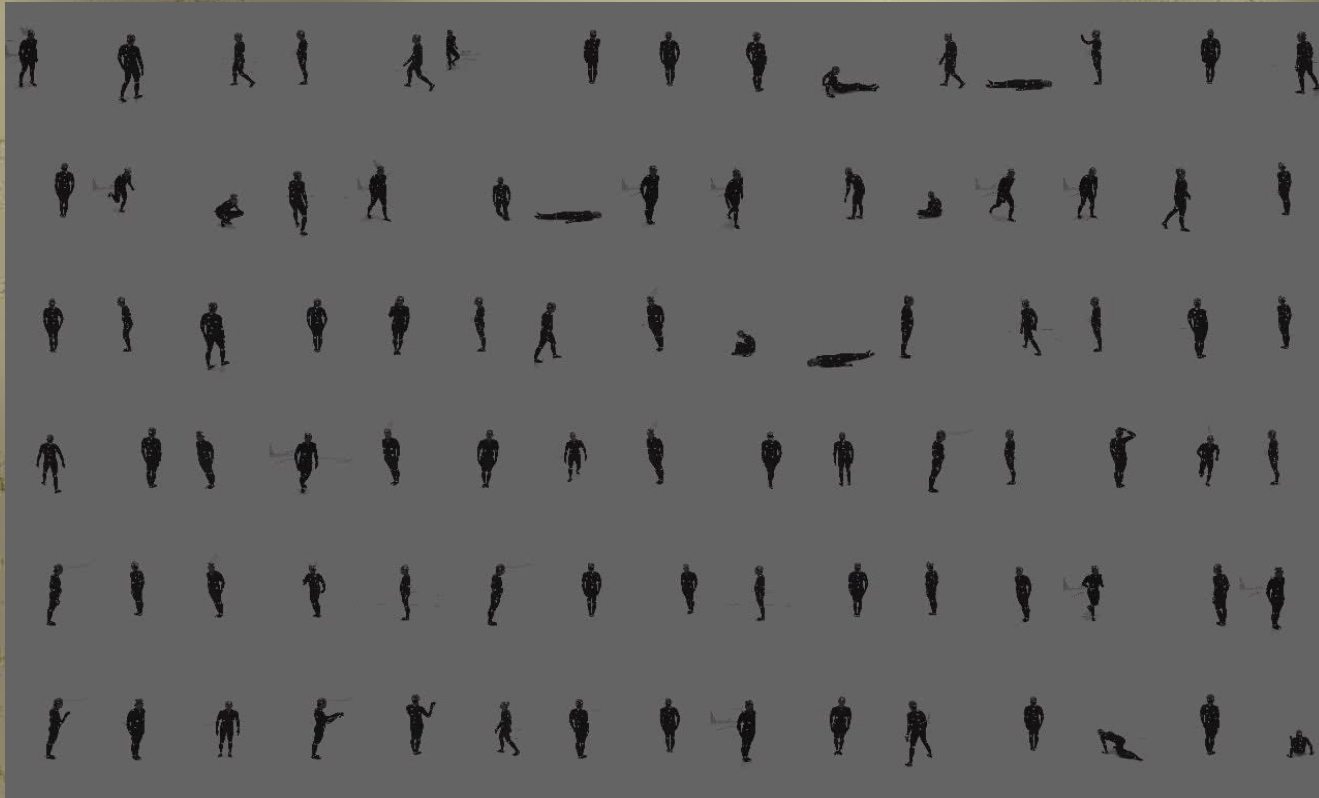
- Easier to solve the visual action problem by going first through the motor action problem.
- Human Activity Language (HAL): a new language for human activity.

Spaces for Human Action



Hyperempiricism in Computer Science

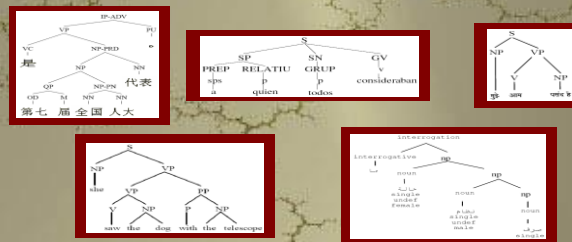
Problem



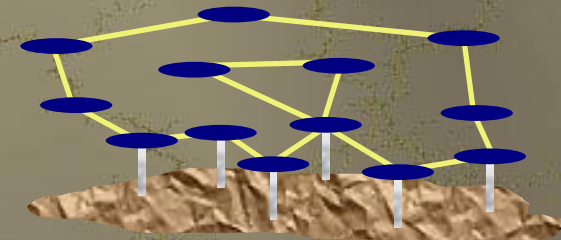
Language Origin



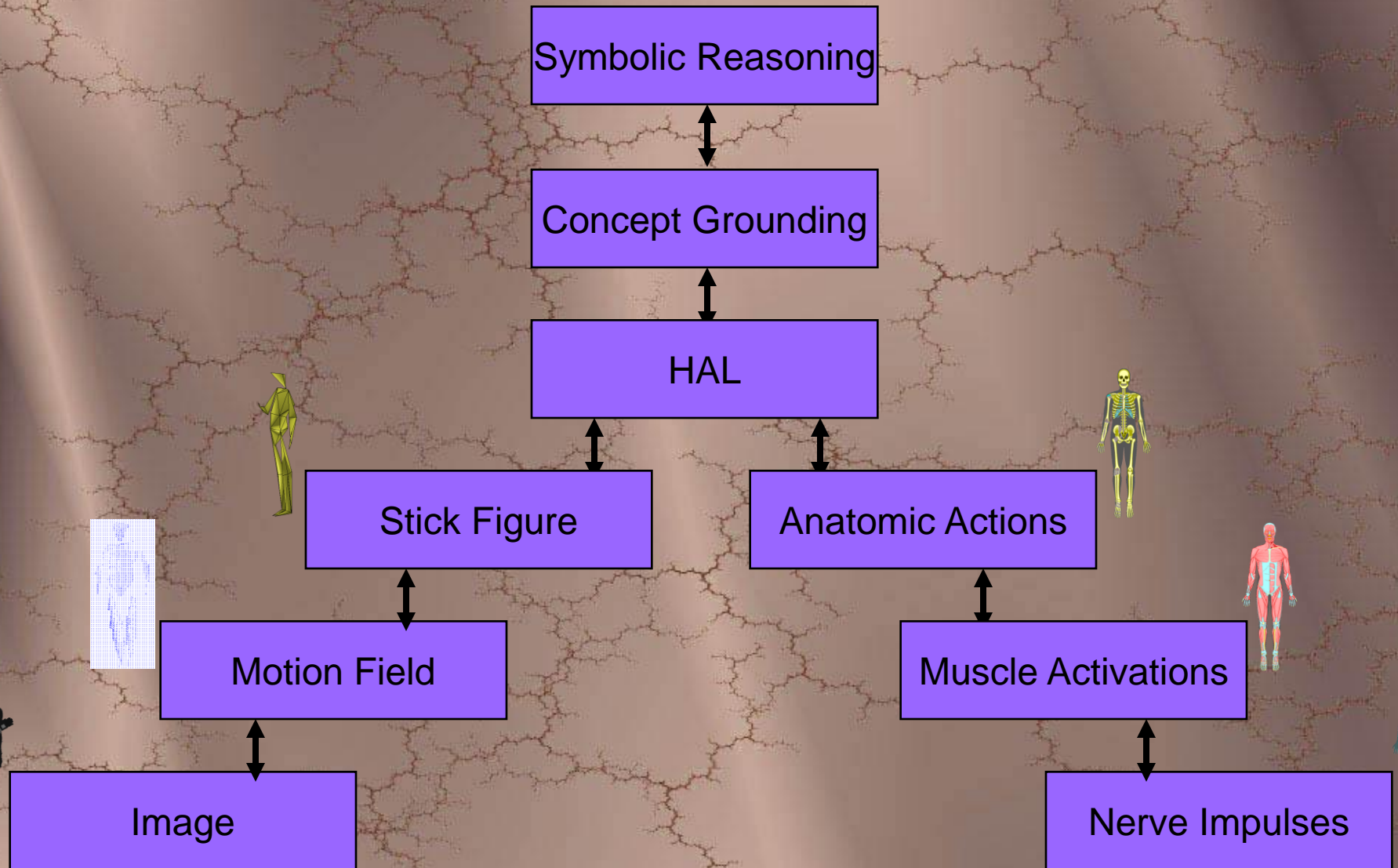
Universal Grammar



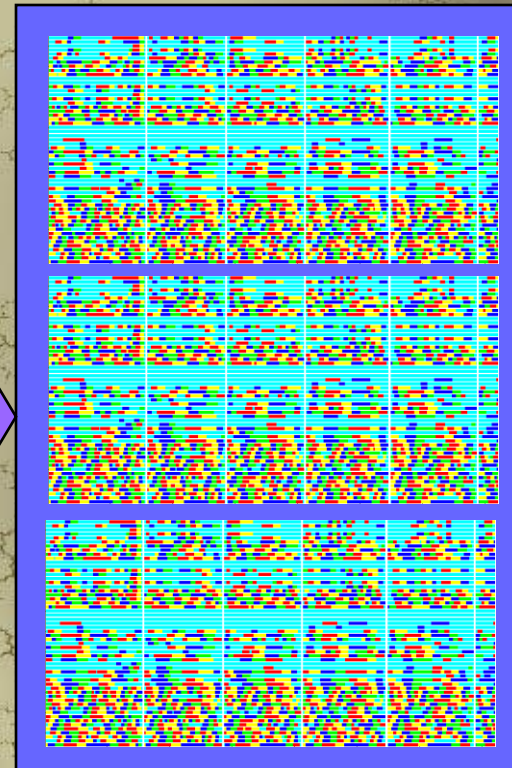
Concept Grounding



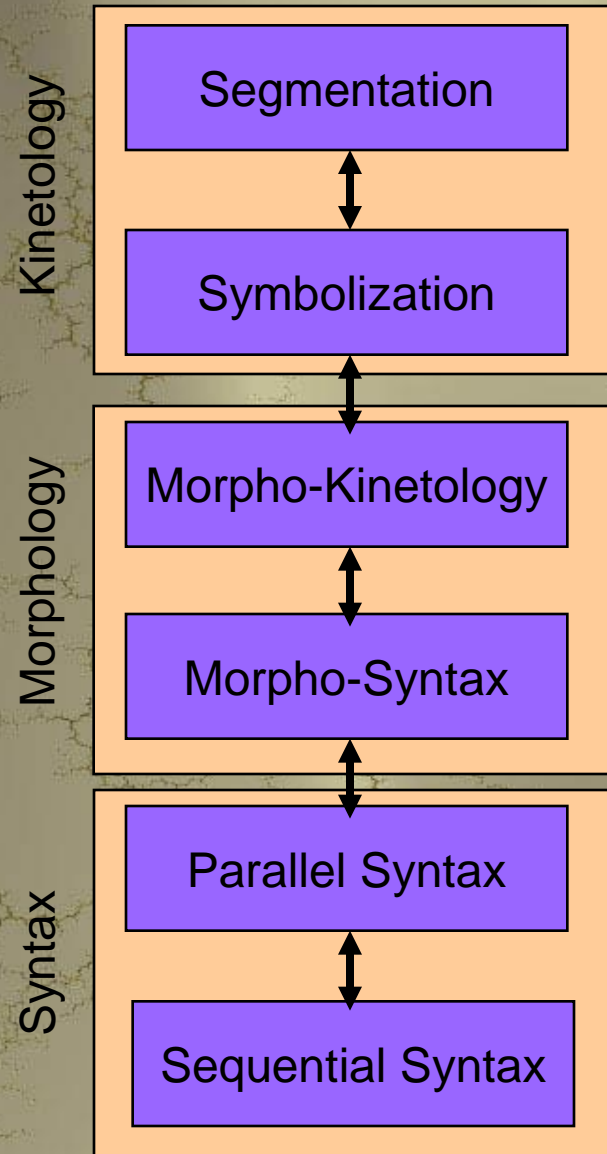
Sensory-Motor Intelligence



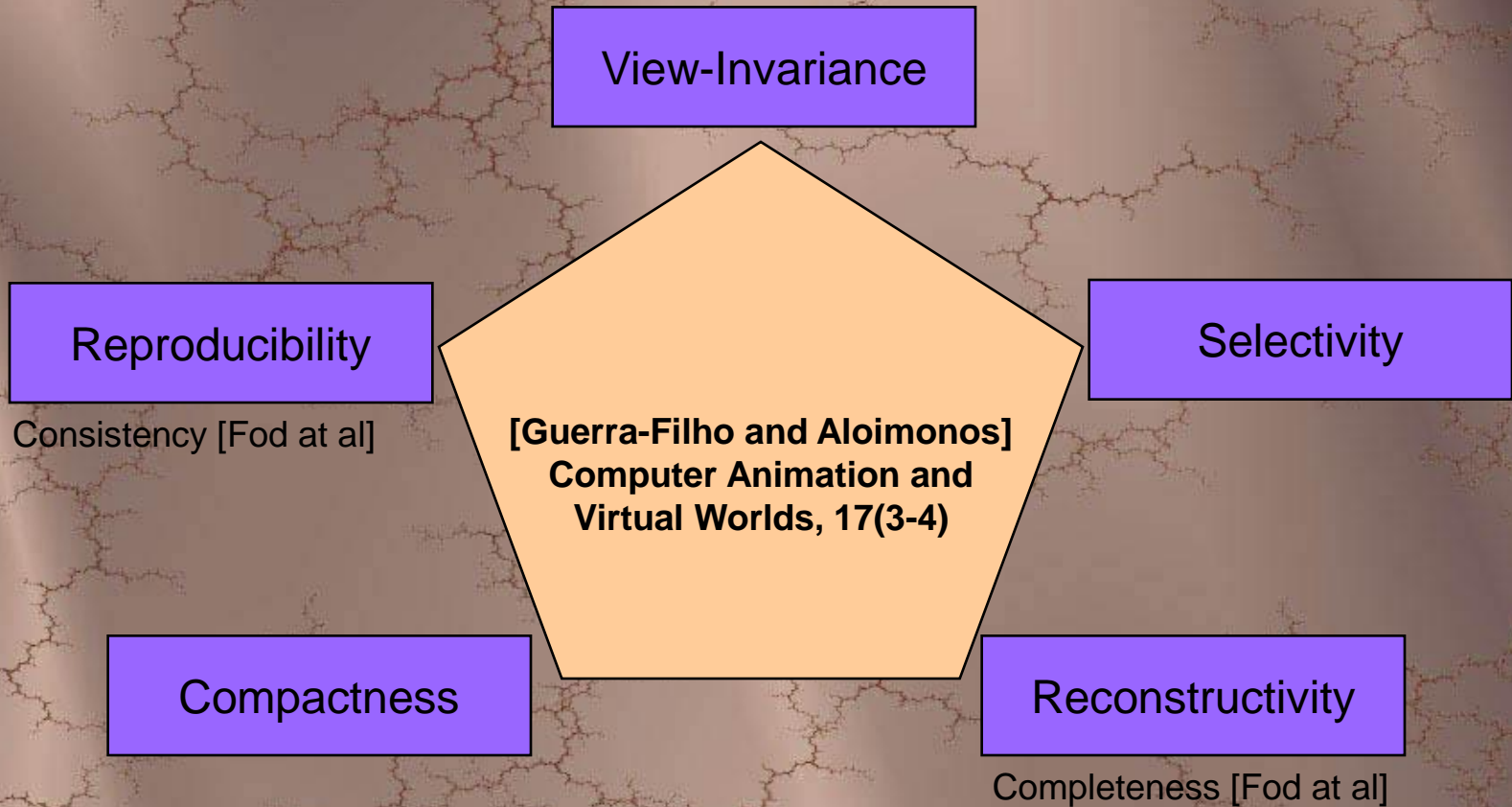
Praxicon



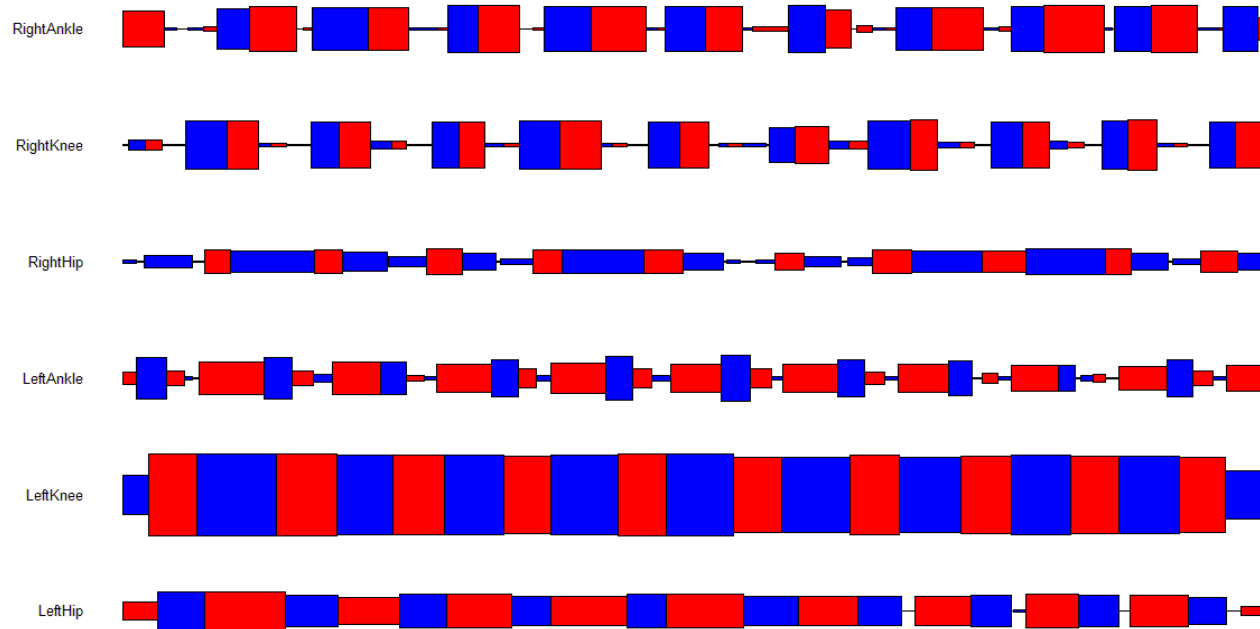
Human Activity Language



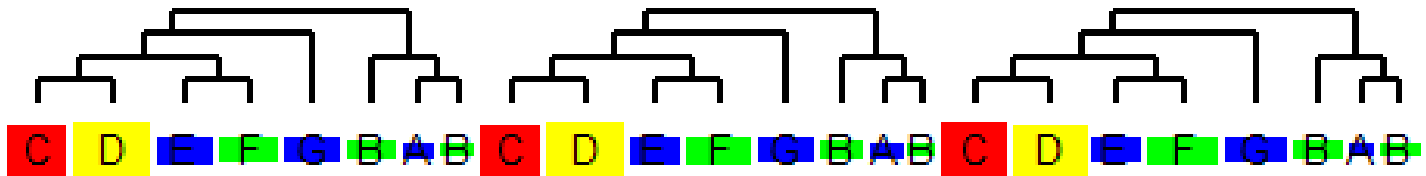
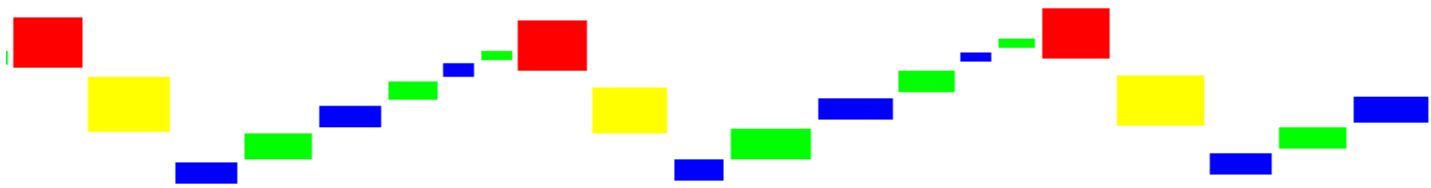
Kinetology



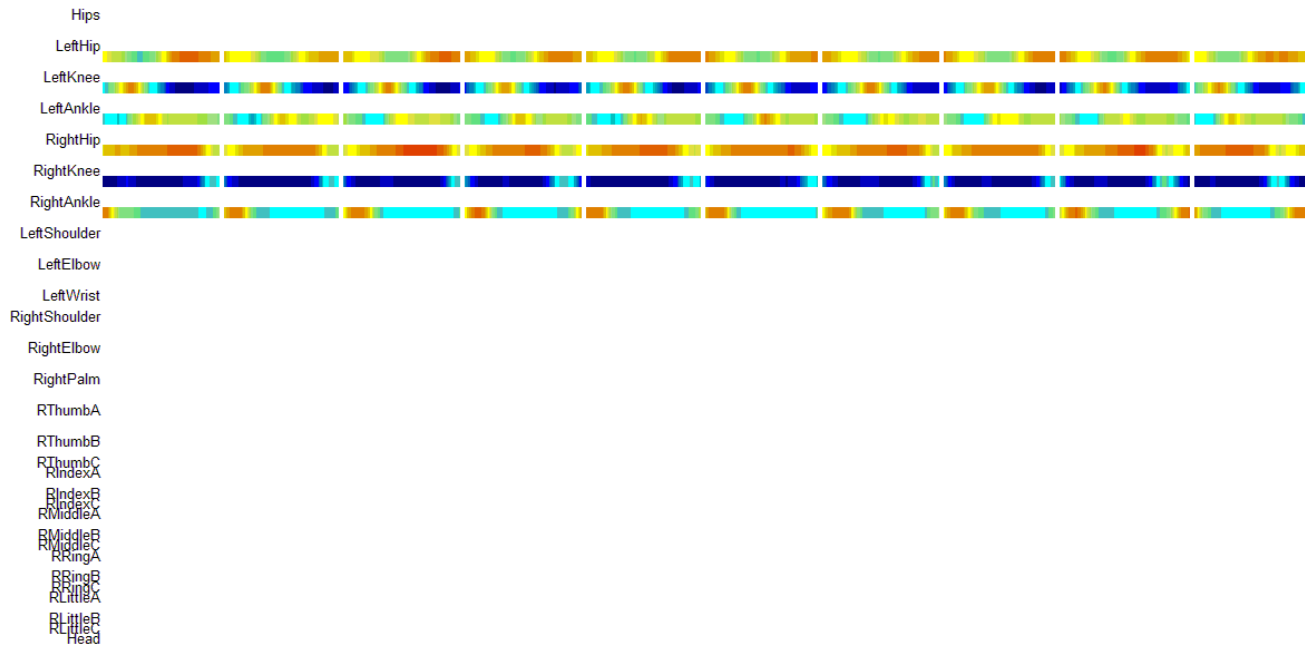
Segmentation



Symbolization

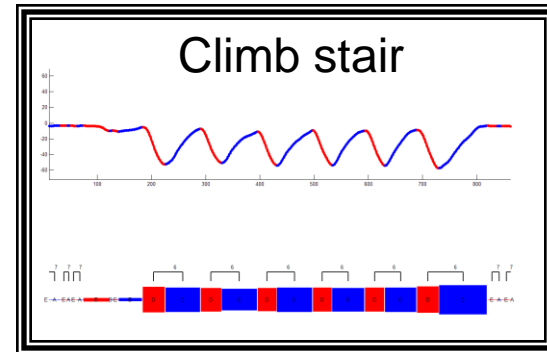
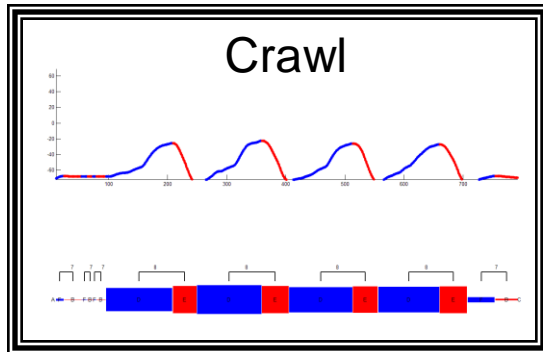
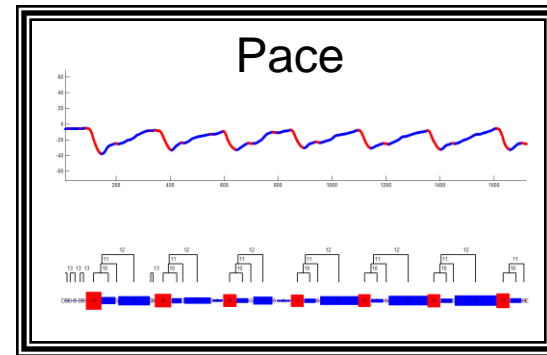
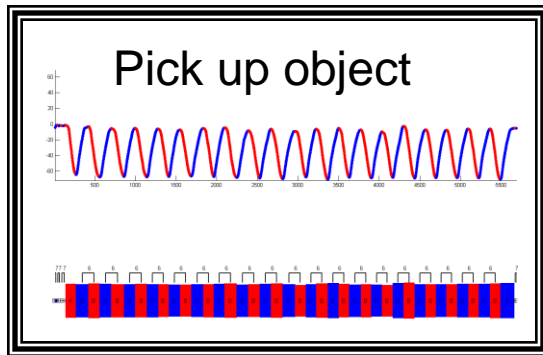


Morphology

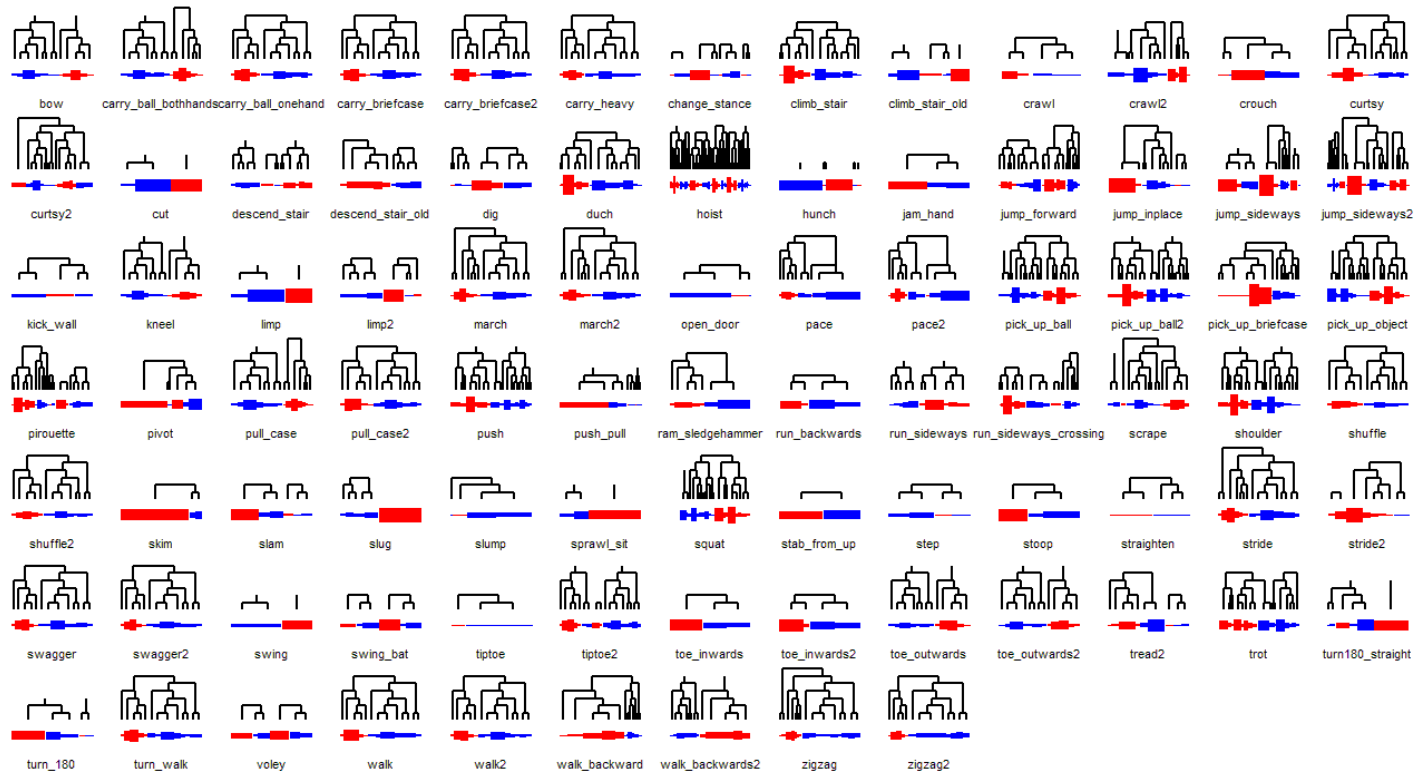


Morpho-kinetology

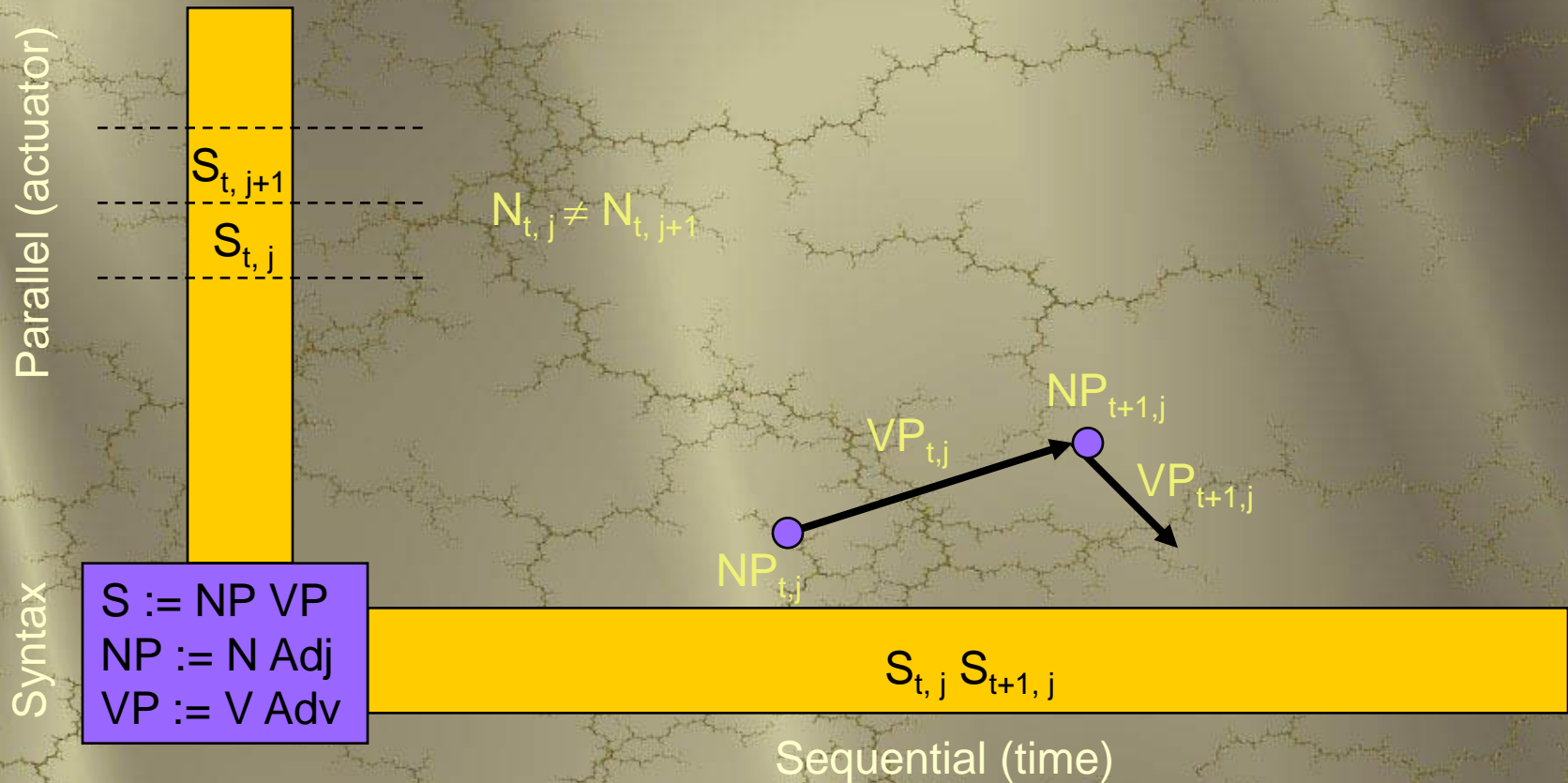
Right Hip Flexion-Extension



Morpho-syntax



Syntax

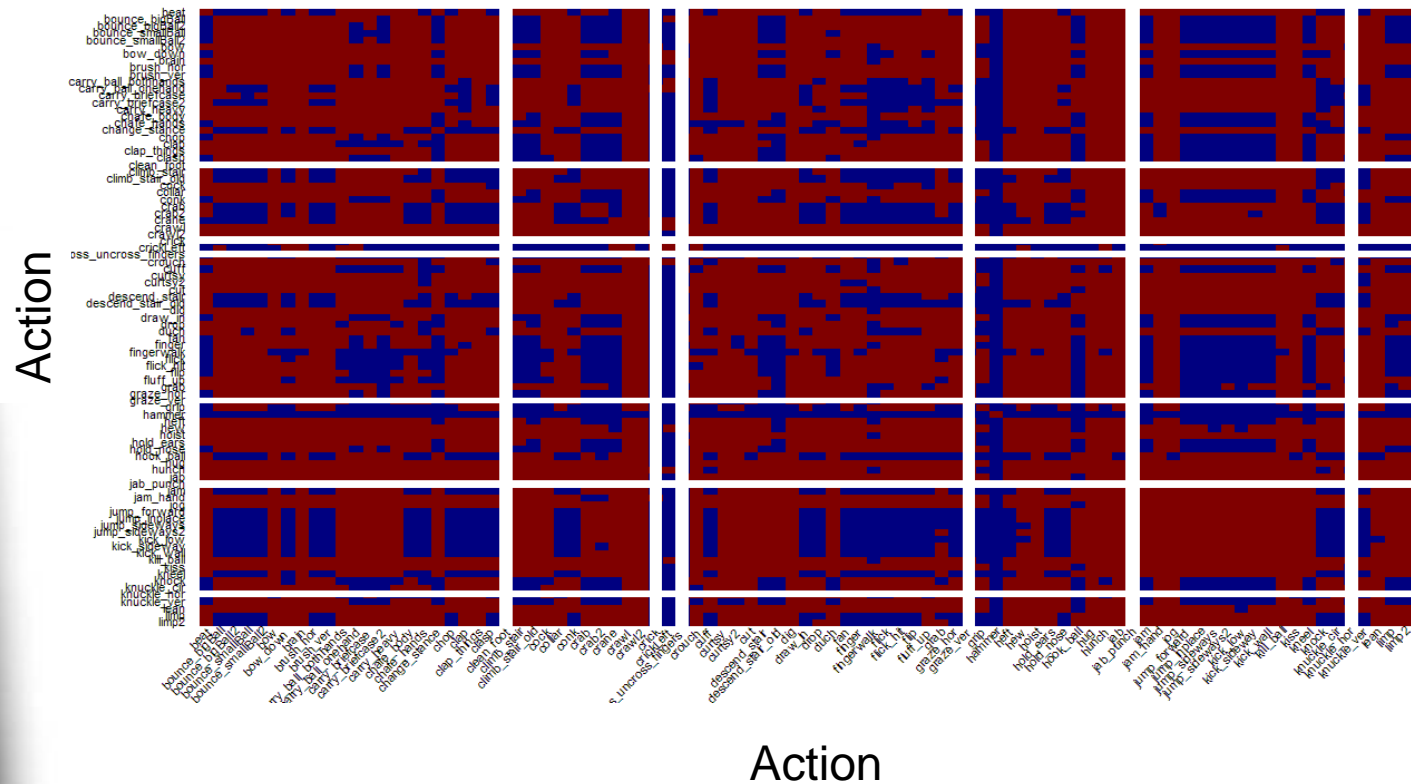


- **Noun**: Body parts active during the execution of a human activity
- **Verb**: Changes each active joint experiences during the activity execution
- **Adjective**: Specifies the initial state of the active joints (initial posture)
- **Adverb**: Modifies verb with purpose of generalization

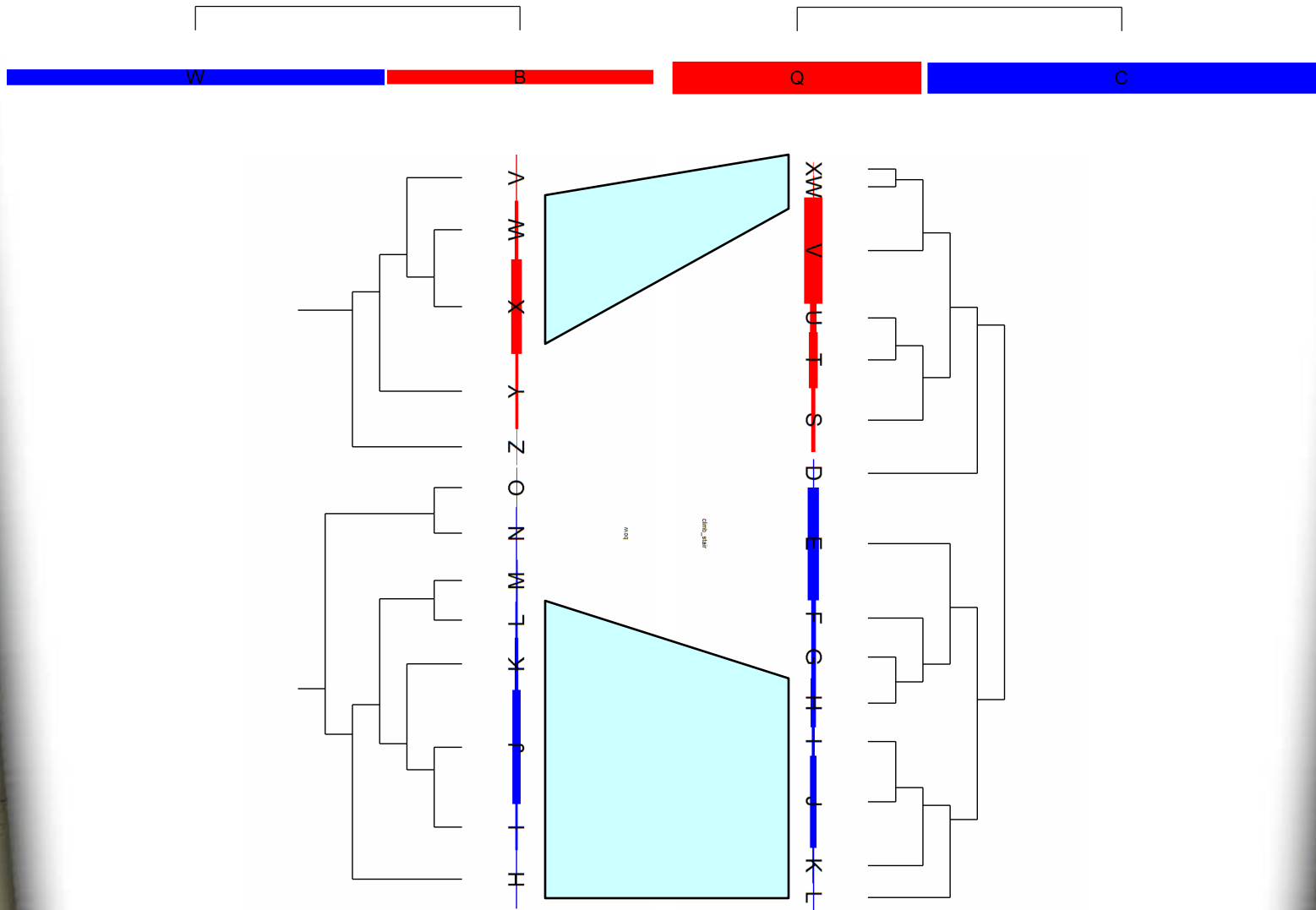
Parallel Syntax

{crick, cross fingers, knuckle, graze, jab, clean foot}

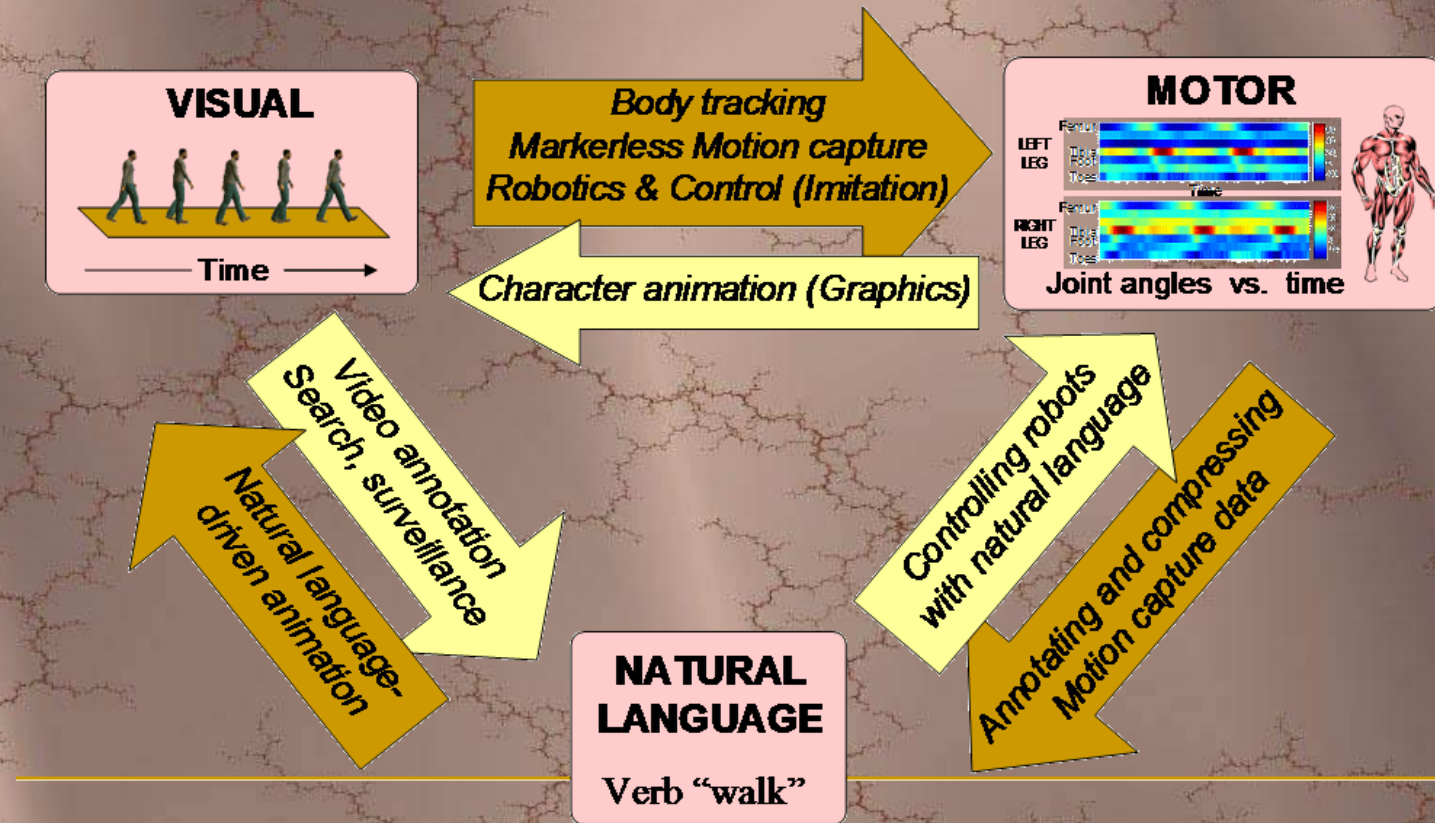
Constraint Matrix



Sequential Syntax



Conclusions



Sensory-Motor Theories vs Symbolic Theories

The Behaviorome Project

