

Gradient Symbol Processing for Phonological Production

or

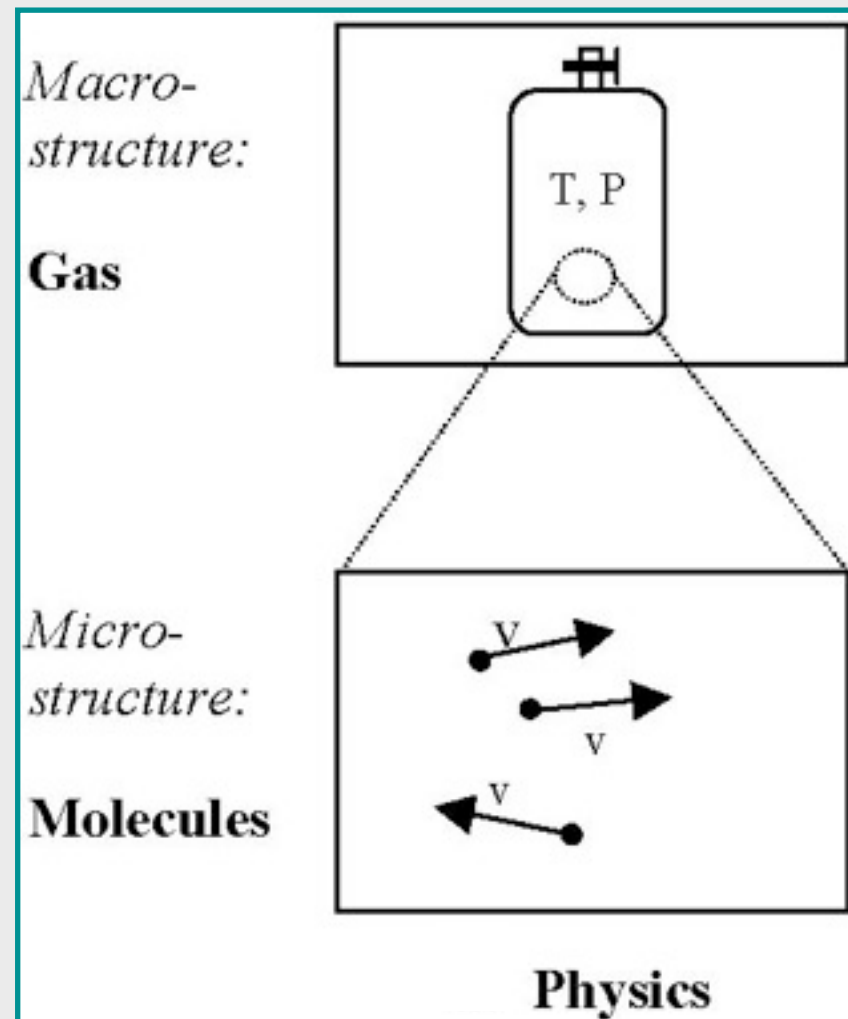
YACA:
Yet Another Cognitive Architecture

Joint work with:

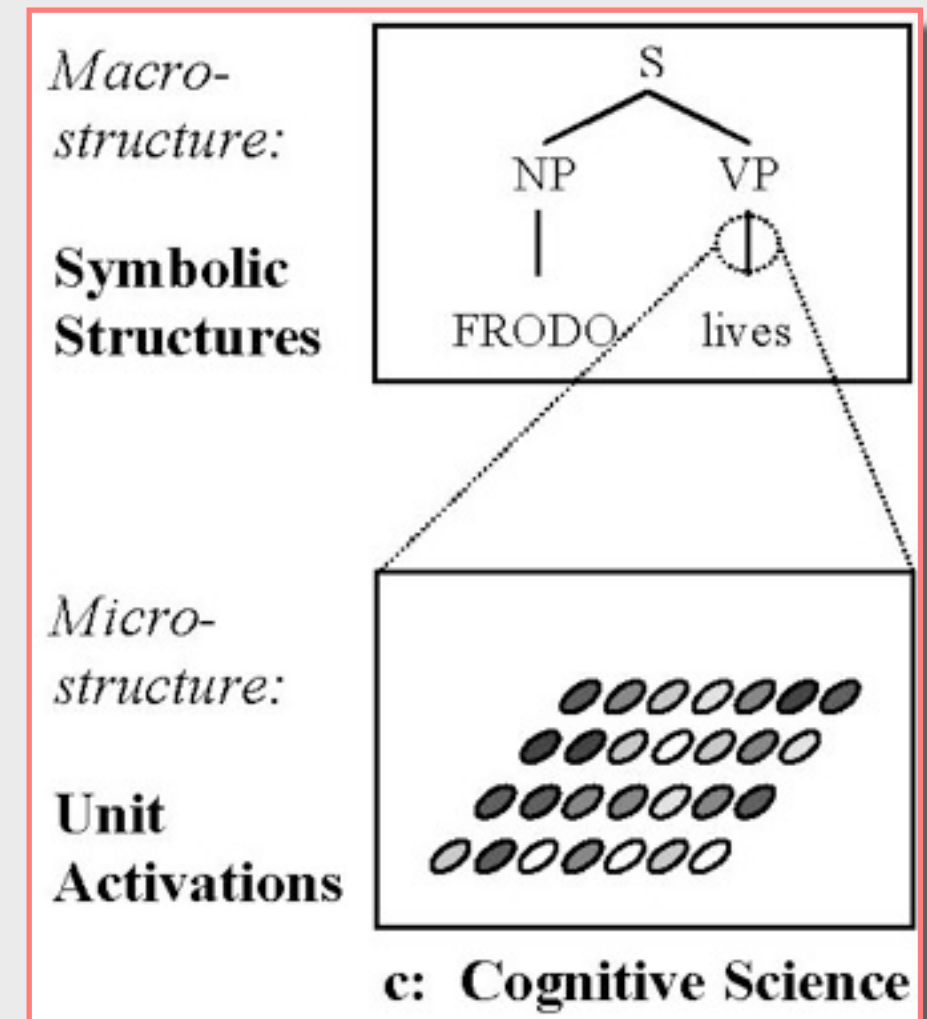
Matt Goldrick (Northwestern Linguistics)
Don Mathis (Johns Hopkins Cognitive Science)

Split-level architecture

The inspiration



The proposal



General cognitive macro-architecture

Graph:

Node:

Representation

information of a particular type
result of function computation
(Markedness constraints)

Edge:

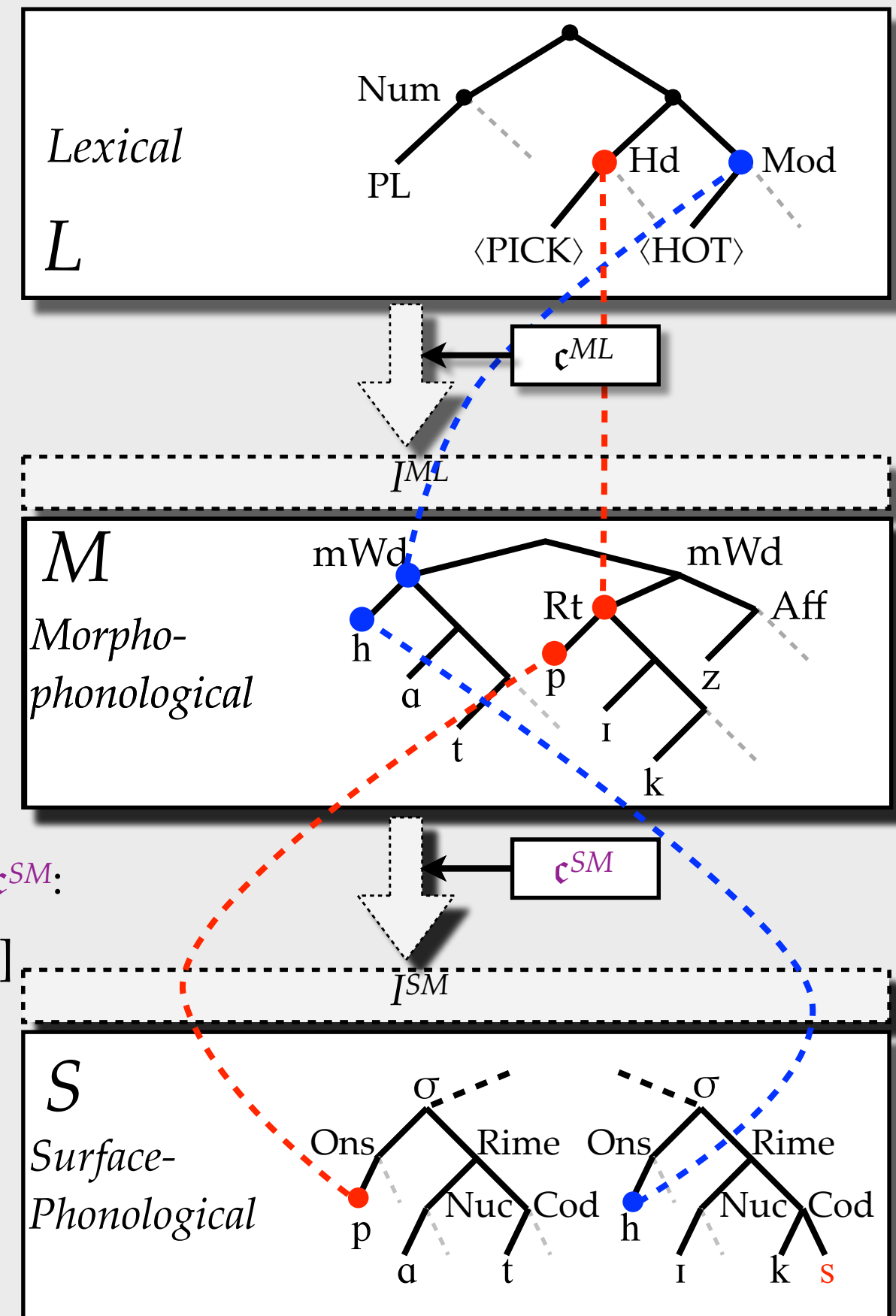
Input to function

Bears a correspondence relation
(Faithfulness constraints)

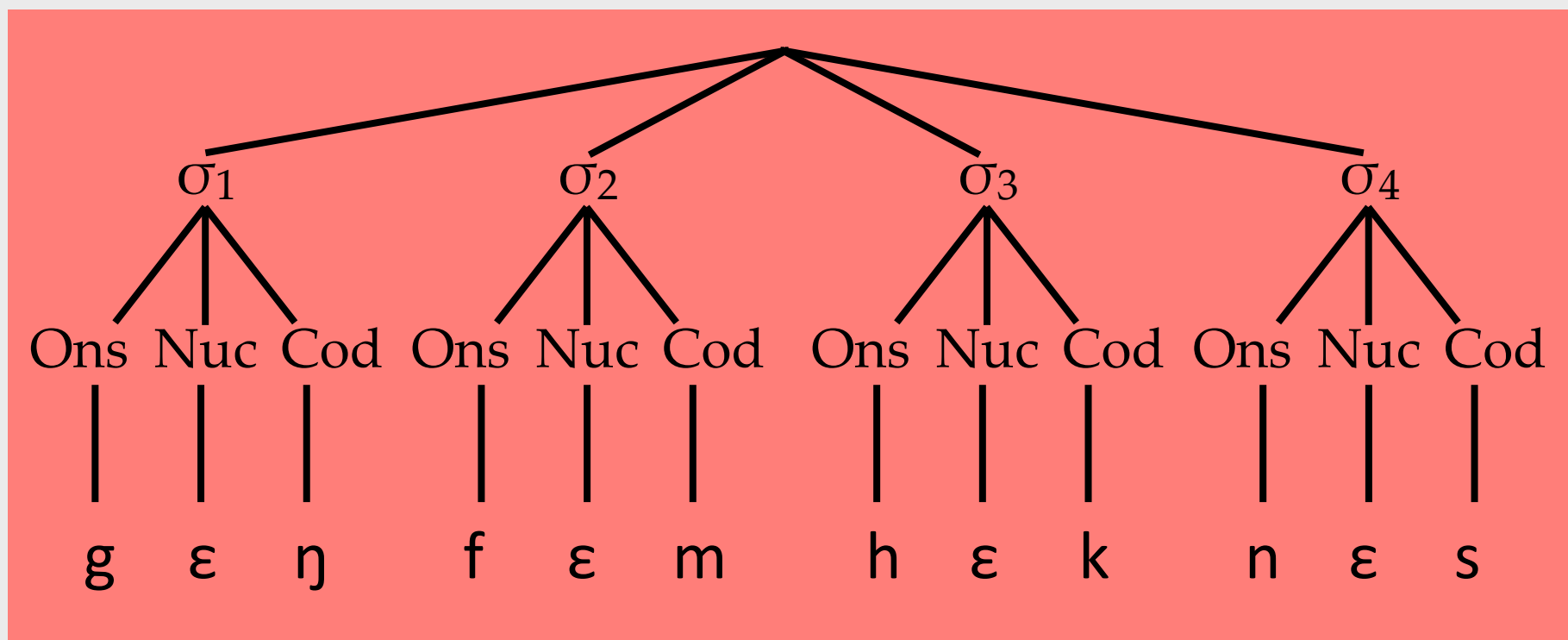
Processes at node:

Optimization

Quantization



Representation



Symbol structures \mathcal{S}

Filler/role decomposition (possibly recursive):

$$s = \{g/[Ons/\sigma_1], \eta/[Cod/\sigma_1], \varepsilon/[Nuc/\sigma_2], \dots\} \subset \mathcal{F} \times \mathcal{R}$$

(activation-)vector space embedding $\mathbf{v}_s \in \mathbf{F} \otimes \mathbf{R} = \mathbb{R}^n$ Random* vectors

$$[g\varepsilon\eta \text{ } f\varepsilon m \dots]: \mathbf{v}_s = \mathbf{g} \otimes \mathbf{r}_{Ons/\sigma_1} + \boldsymbol{\eta} \otimes \mathbf{r}_{Cod/\sigma_1} + \boldsymbol{\varepsilon} \otimes \mathbf{r}_{Nuc/\sigma_1} + \mathbf{f} \otimes \mathbf{r}_{Ons/\sigma_2} + \dots$$

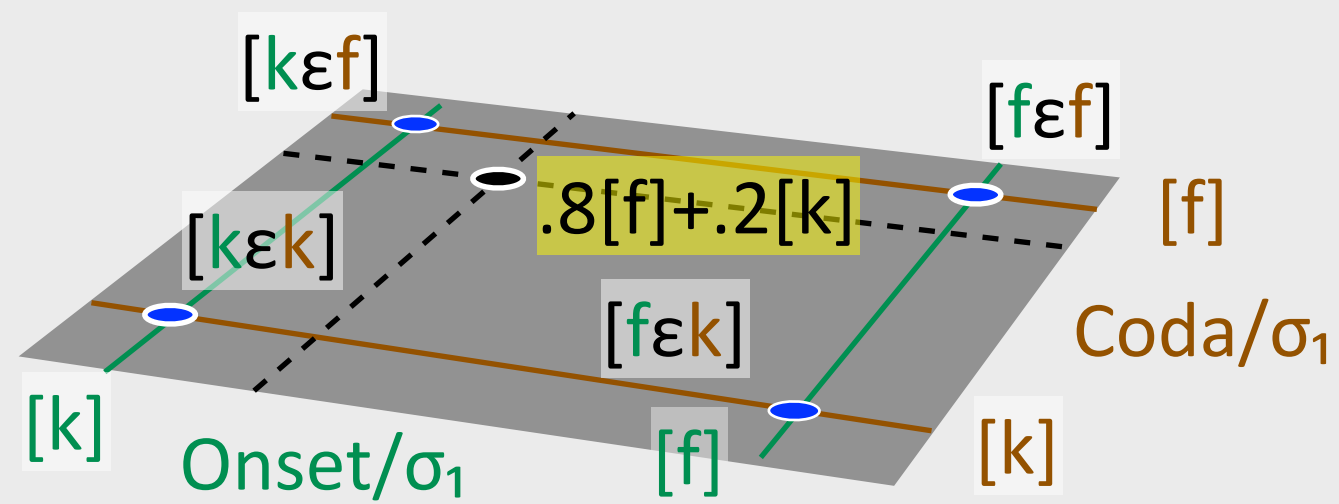
Image of embedding: 'the Grid' # of 'pure states'

$$\mathbf{r}_{Ons/\sigma_2} = \mathbf{r}_{Ons} \otimes \mathbf{r}_{\sigma_2}$$

* Capturing the similarity structure of roles (including recursive hierarchical structure) is a major feature of distributed tensor product representations

Representation

‘the Grid’



Gen: representations

Con: (OT grammar $\mathcal{G} \rightarrow$)

① HG grammar $H_{\mathcal{G}}$

② $\mathcal{G} \rightarrow \mathbf{W}_{\mathcal{G}}$ weight matrix of a network \mathcal{N} s.t.

$$H_{\mathcal{N}_0}(\mathbf{s}) = H_{\mathcal{G}}(s) \quad \forall \mathbf{s} \in \# \text{ (= the Grid) } \text{ — } \textit{iso-Harmonic embedding}$$

Theorem. For any deterministic neural network in a certain class, during processing, Harmony continuously increases, reaching a *local* optimum.

- This is **network Harmony** $H_{\mathcal{N}} = H_{\mathcal{N}_0} + H_{\mathcal{N}_1}$ $H_{\mathcal{N}_1}(\mathbf{a}) = \frac{1}{2}|\mathbf{a}|^2$

$$H_{\mathcal{N}_0}(\mathbf{a}) \equiv \sum_{\beta\gamma} a_{\beta} W_{\beta\gamma} a_{\gamma} \quad \text{— quadratic, dependent on } \mathbf{W}$$

Theorem. For any stochastic neural network in a certain class, during processing, the probability of visiting a state \mathbf{a} approaches

$$p(\mathbf{a}) \propto e^{H_{\mathcal{N}}(\mathbf{a})/T} \quad (T = \text{randomness parameter})$$

As $T \rightarrow 0$, the probability the network is in a *globally* optimal state $\rightarrow 1$.

- These networks use a **Diffusion Dynamics**

$$da_{\beta} = \sum_{\gamma} W_{\beta\gamma} a_{\gamma} dt + \sqrt{2T} dB_{\beta} = \frac{\partial H_{\mathcal{N}}}{\partial a_{\beta}} dt + \sqrt{2T} dB_{\beta}$$

Processing: **Diffusion Dynamics**

- State moves in time so as to increase Harmony $H(s)$, on average
 - ♦ randomness in state changes with variance $\propto T$
 - ♦ during processing, $T \rightarrow 0$
 - ♦ hence $p(s) \rightarrow 0$ except for the state(s) with maximal Harmony
 - ♦ N.B.: randomness needed to find *global* Harmony maxima
 - ♦ not infallible: errors occur
 - ♦ from mechanism responsible for *correct* performance

 **Optimization** process

Theorem. Any rewrite-rule grammar can be expressed as a second-order Harmonic Grammar.

Theorem. For any second-order Harmonic Grammar $H_{\mathcal{G}}$, we can construct a recurrent network \mathcal{N} with a harmony function $H_{\mathcal{N}}$ that provides an iso-Harmonic embedding

i.e., yields the same values as $H_{\mathcal{G}}$ on every pure (grid) state s :

$$H_{\mathcal{N}}(\mathbf{s}) = H_{\mathcal{G}}(s)$$

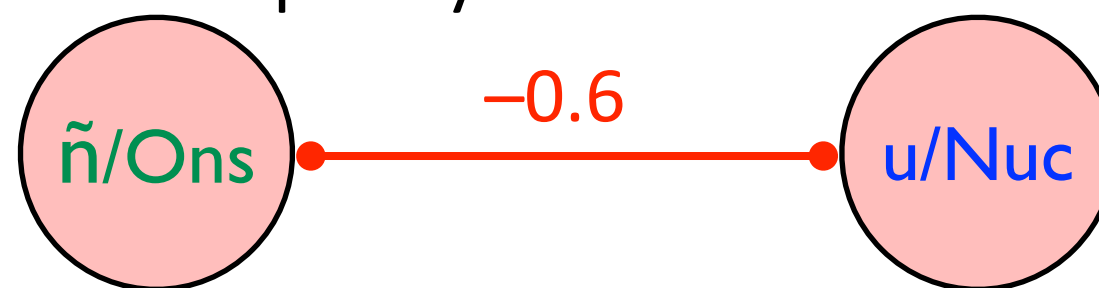
Corollary. For any Harmonic Grammar $H_{\mathcal{G}}$, we can construct a recurrent network \mathcal{N} such that as $T \rightarrow 0$, the probability the network is in a gradient state that is *globally* optimal w.r.t. $H_{\mathcal{G}} \rightarrow 1$.

Spreading activation = finding optimal solution to weighted constraints

E.g., phonotactic constraint: * $\tilde{n}u$ (American *muse* vs. *news*)

Harmony maximization as constraint satisfaction

Consider this connection in a purely localist network:



Same constraint with distributed representations:

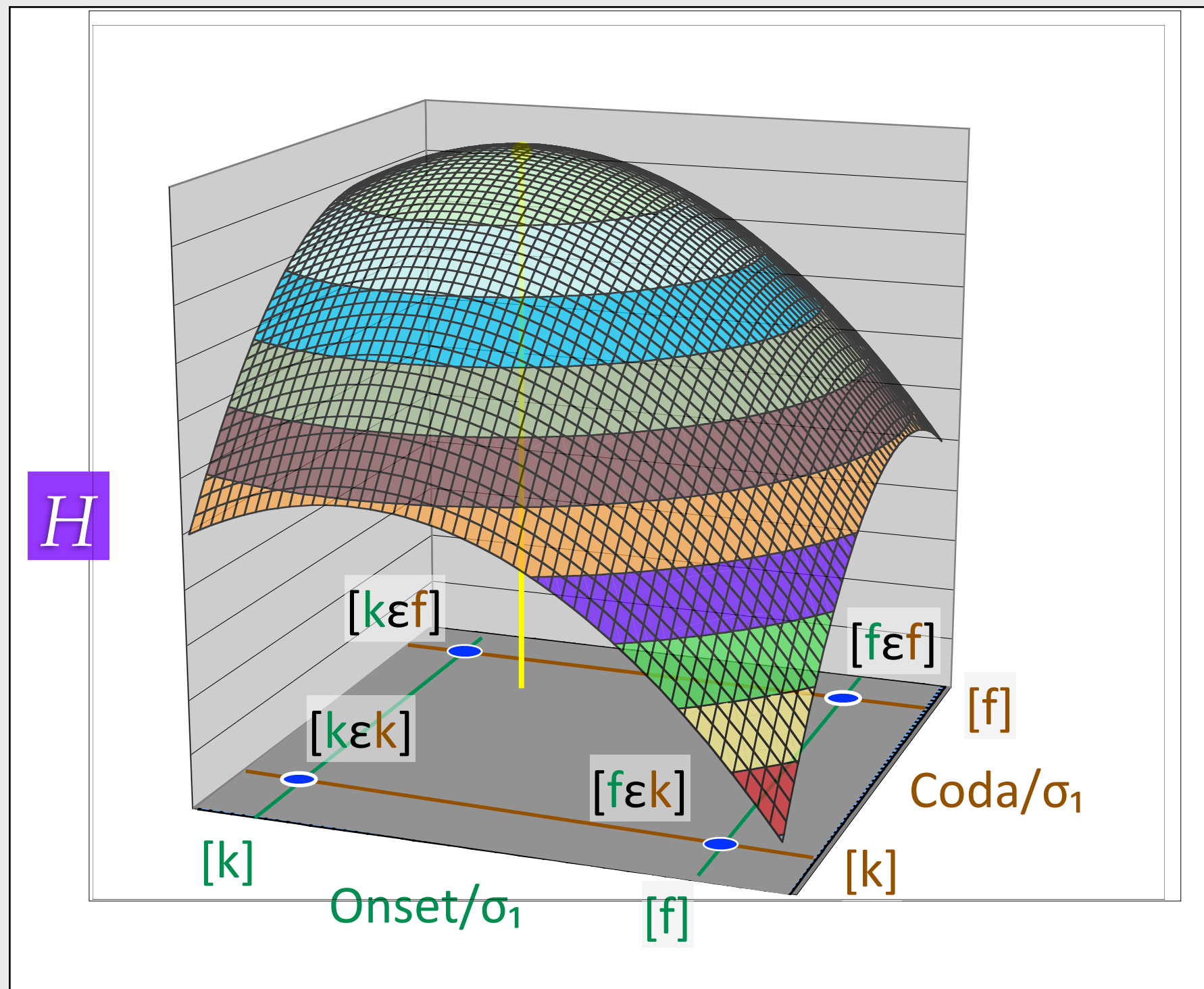
the weight matrix **W** such that

when activation patterns are re-described in a new coordinate system
in which the representations become local,
W becomes equal to the connection above.

Problems:

want a maximum of H at every grammatical structure

but H is quadratic: it can have only one global maximum



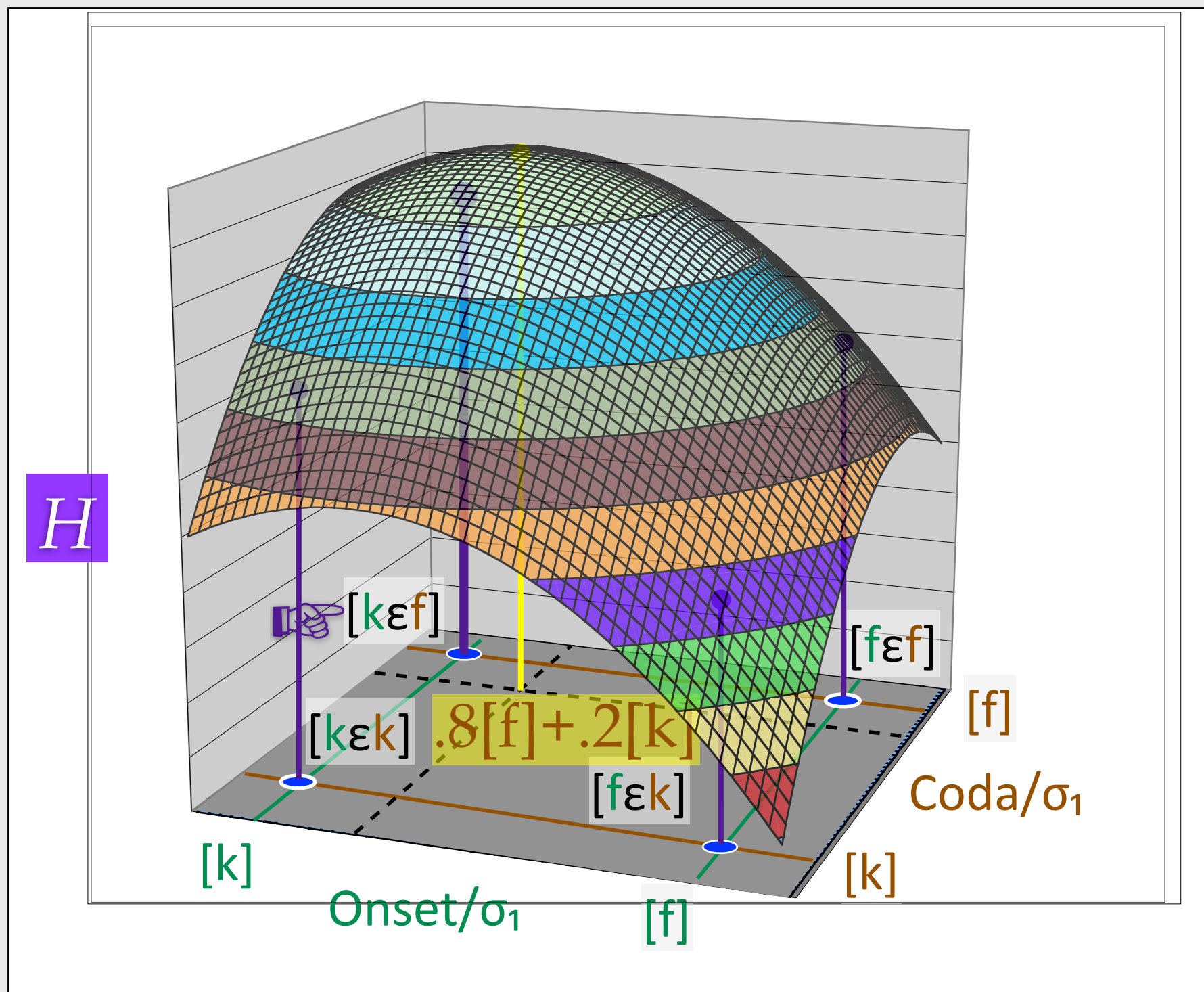
Problems:

want a maximum of H at every grammatical structure

but H is quadratic: it can have only one global maximum

$\operatorname{argmax}_{\mathbf{a} \in \mathbb{R}^n} H_{\mathcal{N}}(\mathbf{a}) \notin \text{Grid}$: it is a *blend* state

H restricted to the grid **can** have multiple maxima



Corollary. For any Harmonic Grammar H_G , we can construct a recurrent network \mathcal{N} such that as $T \rightarrow 0$, the probability the network is in a (gradient) state that is *globally* optimal w.r.t. $H_G \rightarrow 1$.

This is a **blend** of well-formed constituents, not a globally coherent pure state.
(A general problem, not limited to grammars.)

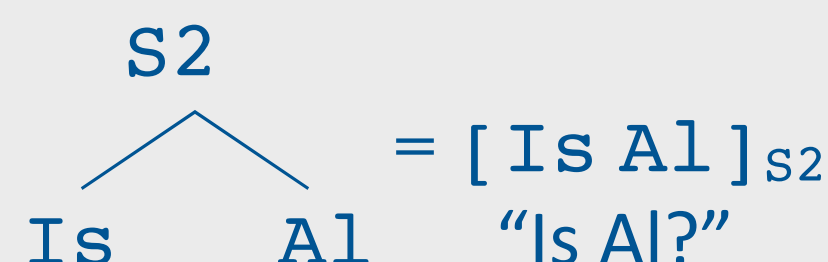
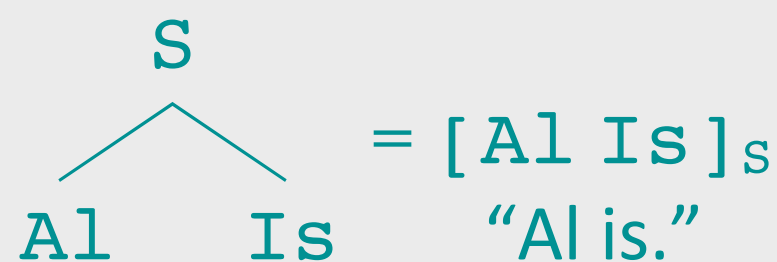
A nanogrammar \mathcal{G}

Its nanolanguage \mathcal{L}

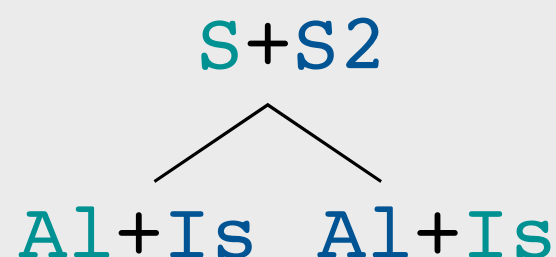
Start symbols: $\{S, S2\}$

$S \rightarrow Al\ Is$

$S2 \rightarrow Is\ Al$



The global H optimum is proportional to



This is why we need **quantization**.

Problems:

want a maximum of H at every grammatical structure

but H is quadratic: it can have only one global maximum

$\operatorname{argmax}_{\mathbf{a} \in \mathbb{R}^n} H_{\mathcal{N}}(\mathbf{a}) \notin \text{Grid}$: it is a *blend* state

H restricted to the grid **can** have multiple maxima

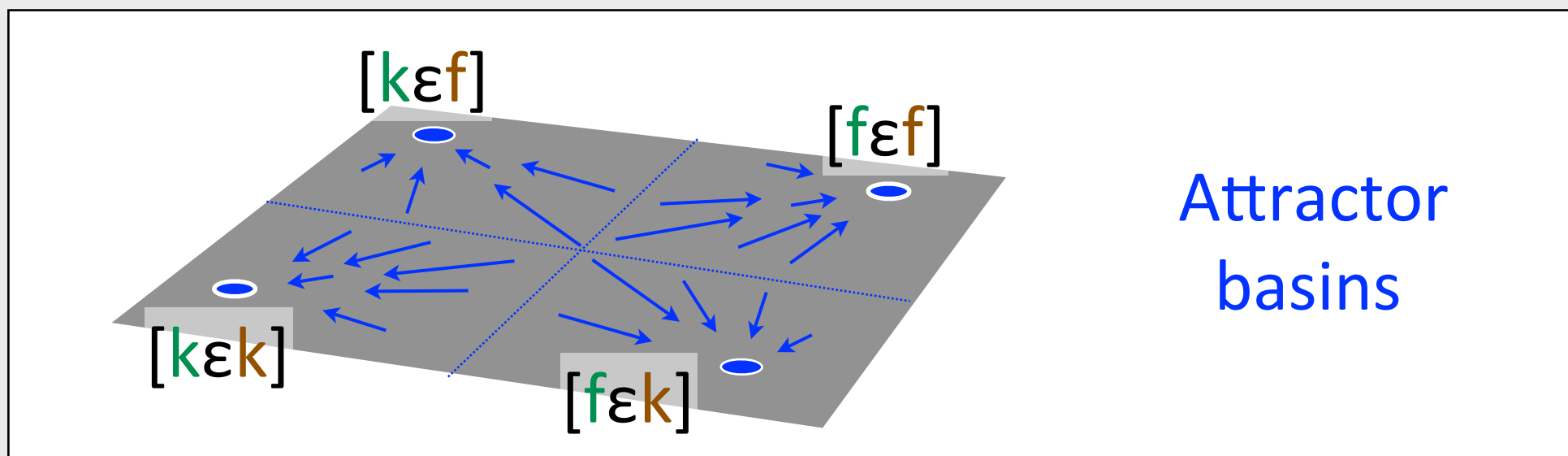
H is meaningful only on the grid

Proposal:

Add a *quantization* dynamics with an attractor at every $s \in \text{Grid}$

Discretization Dynamics

- A spreading activation algorithm that creates an **attractor** at all and only the points of the grid
- Isotropic/symmetric/all attractors equivalent:
 - ✦ *optimization dynamics* pushes towards correct basin



- Distributed winner-take-all network (non-linear mutual inhibition)
 - ✦ Lotka-Volterra equations (Baird & Eeckmann 1993)

$$\frac{dx_{\beta}}{dt} = x_{\beta} - \sum_{\mu\nu} W_{\beta\mu\nu} x_{\mu} x_{\nu} \quad W_{\beta\mu\nu} = \sum_{jk} M_{\beta k} M_{k\mu}^{-1} M_{j\nu}^{-1} (2 - \delta_{jk})$$

$\mathbf{M} = \mathbf{F} \otimes \mathbf{R}$, \mathbf{F} = matrix of symbol (filler) patterns, \mathbf{R} = of position (role)

Harmony Optimization Dynamics

- Diffusion; as processing proceeds, $T \rightarrow 0$
- Pushes towards best gradient (blend) state
 - ✦ ignores discreteness

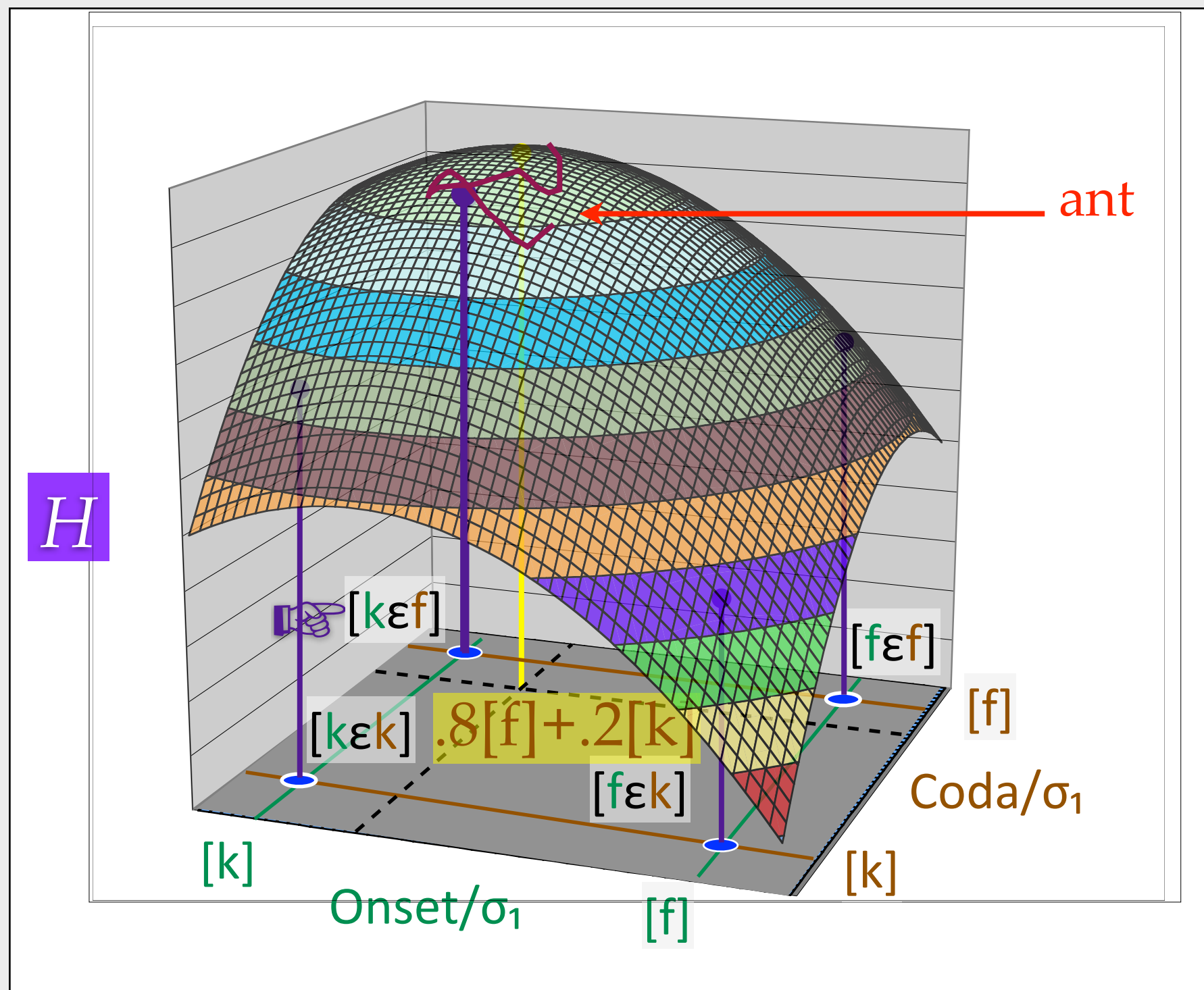
Quantization Dynamics

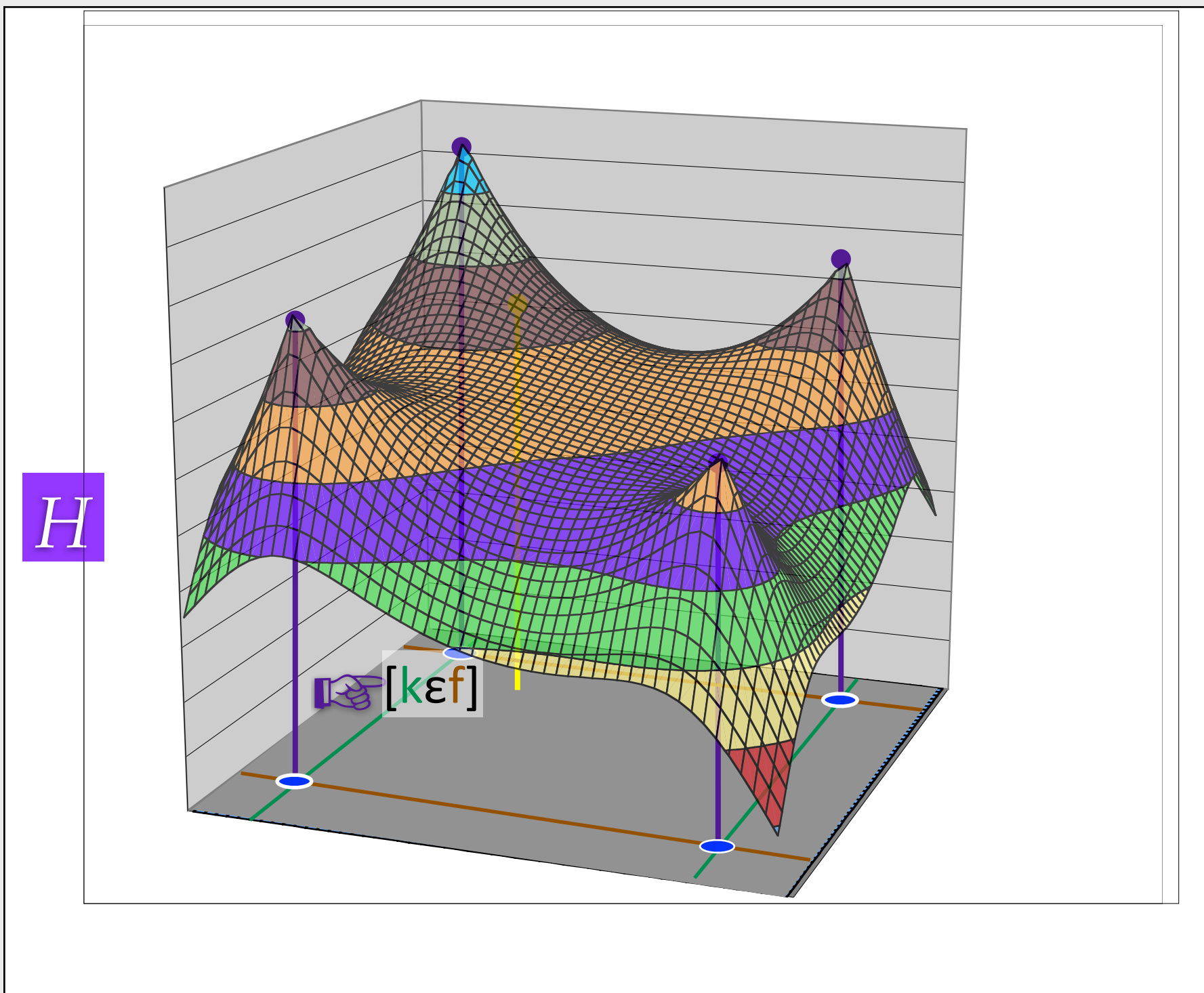
- Pushes towards the grid of discrete (pure) states
 - ✦ ignores well-formedness

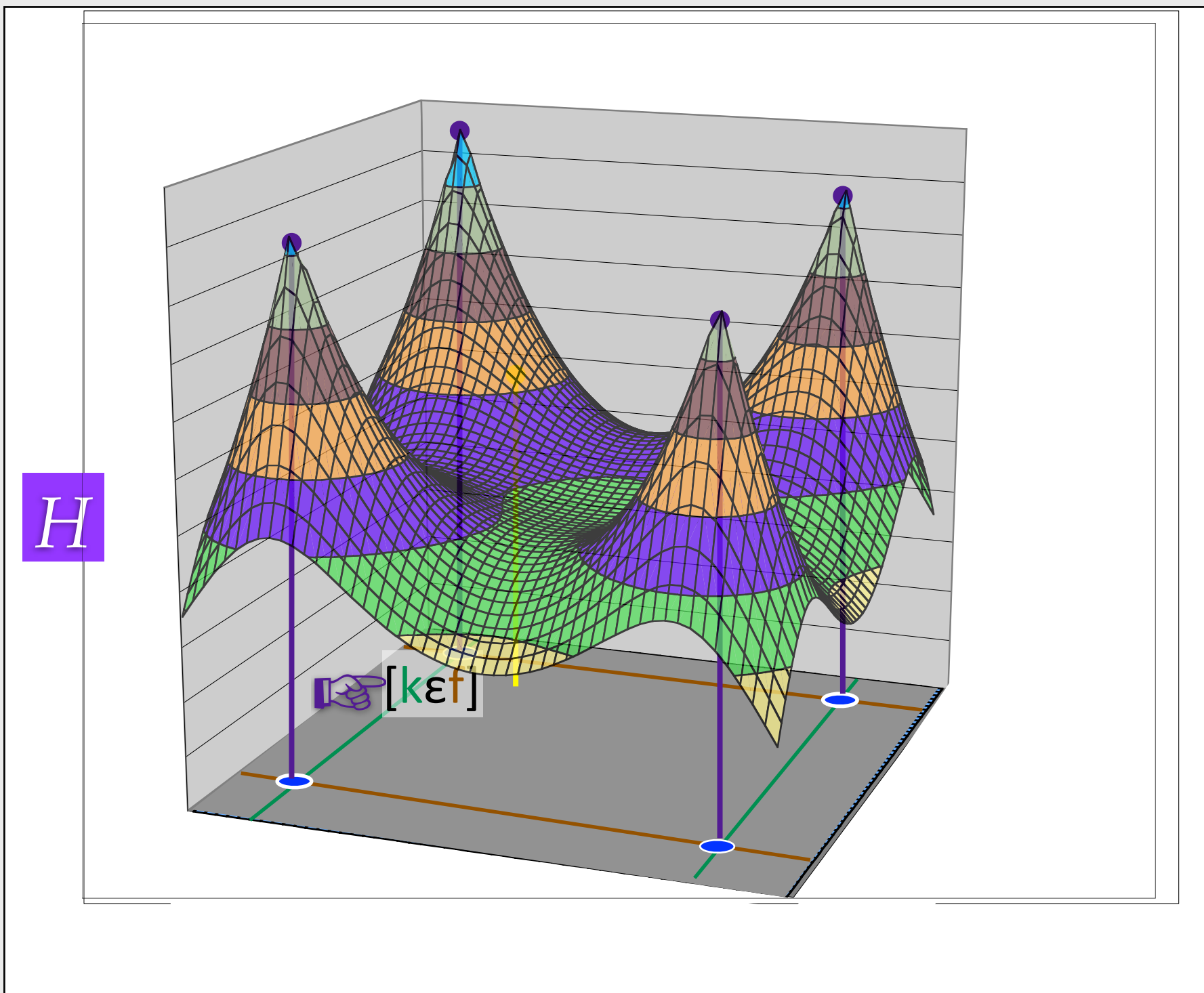
Combination

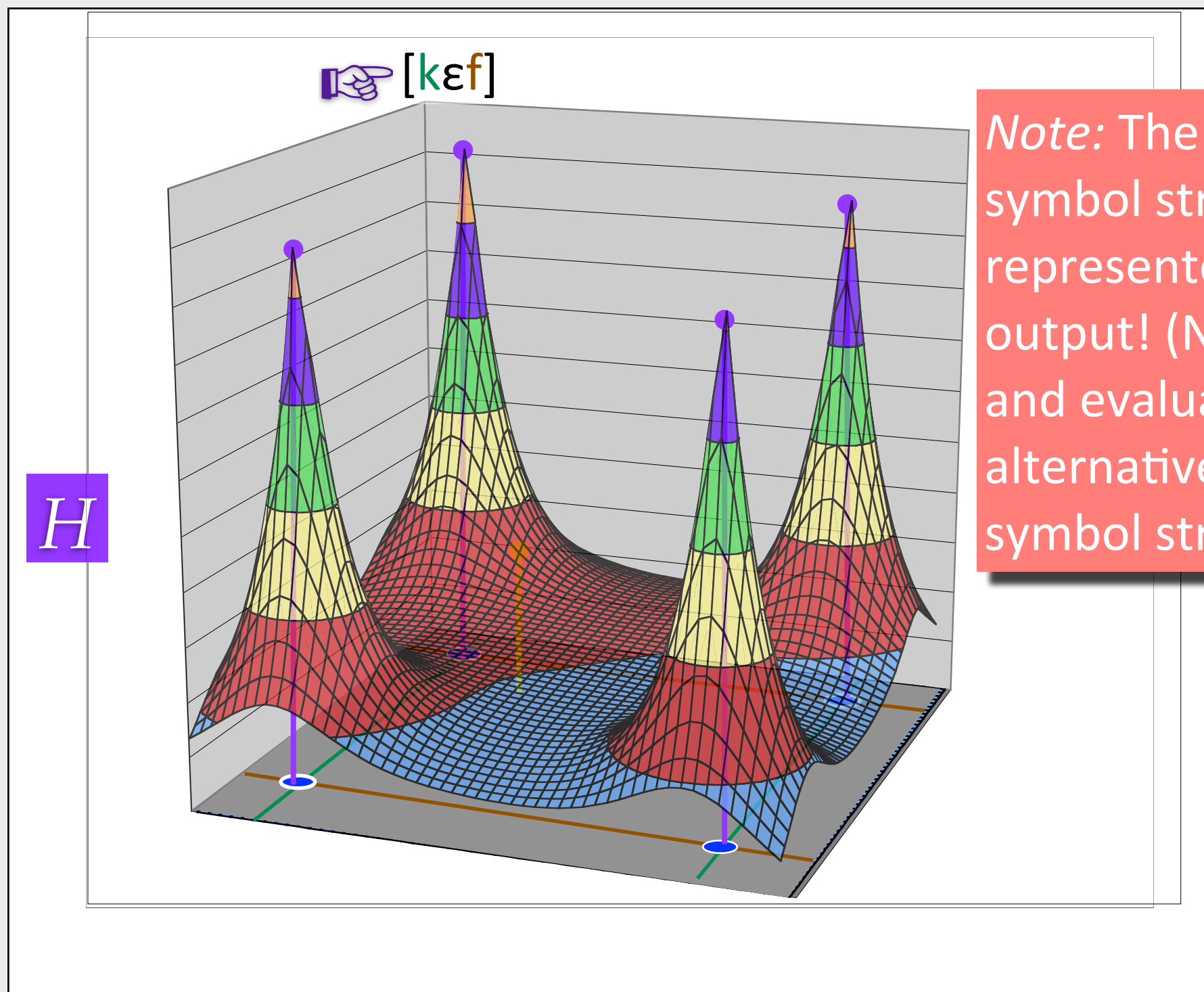
- The weighted sum of these two dynamics
- As processing proceeds, the relative weight of optimization
 - ✦ $\lambda \rightarrow 0$
 - ✦ discretization pressure grows, dominates final computation

$$\mathcal{D} = \lambda \mathcal{D}_{opt} + [1-\lambda] \mathcal{D}_{quant}$$









Note: The only discrete symbol structure ever represented is the final output! (Not: generate and evaluate many alternative discrete symbol structures.)

The ant should end up at the highest peak
— or at an erroneous peak with *prob* $\sim H$

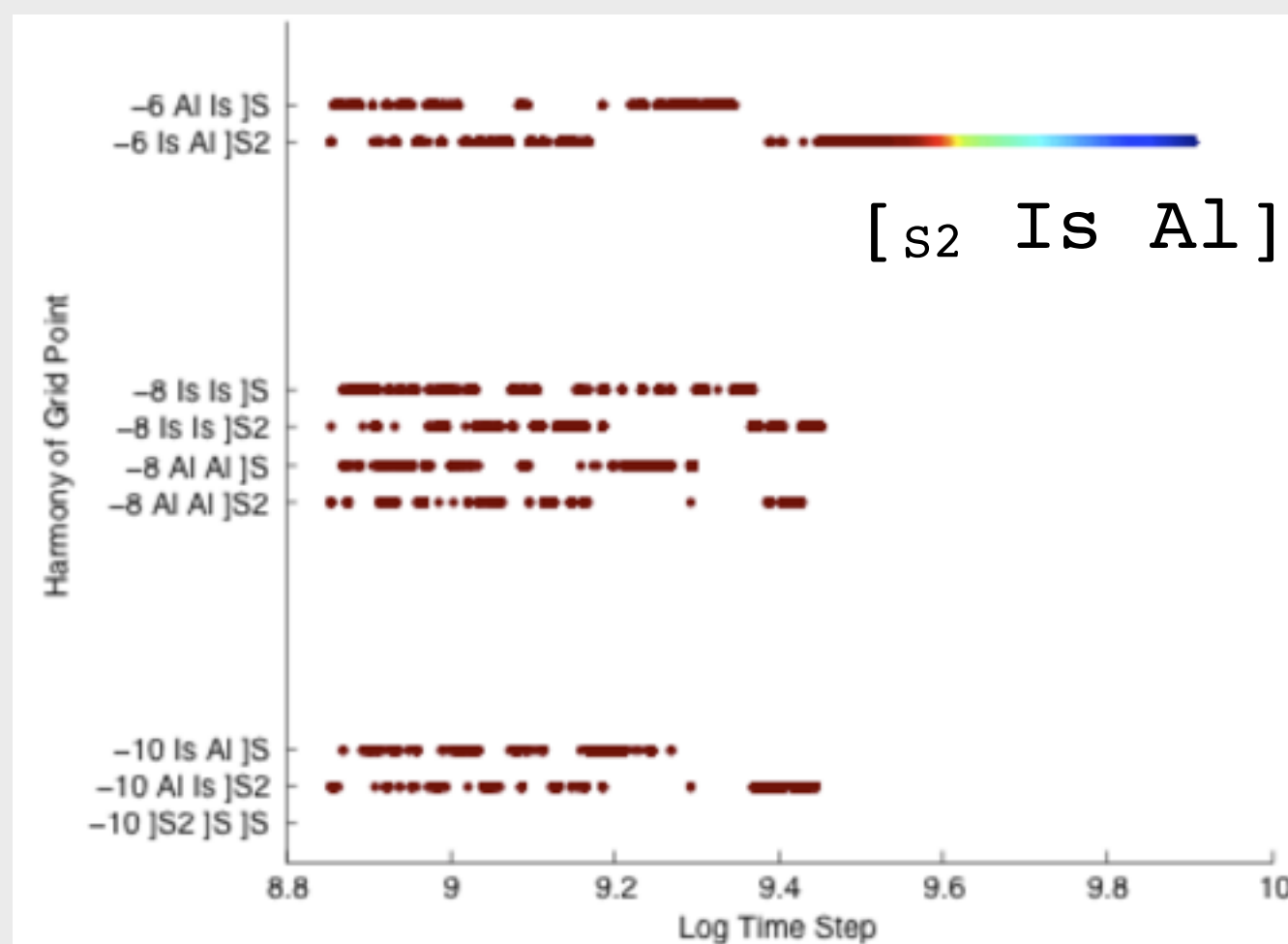
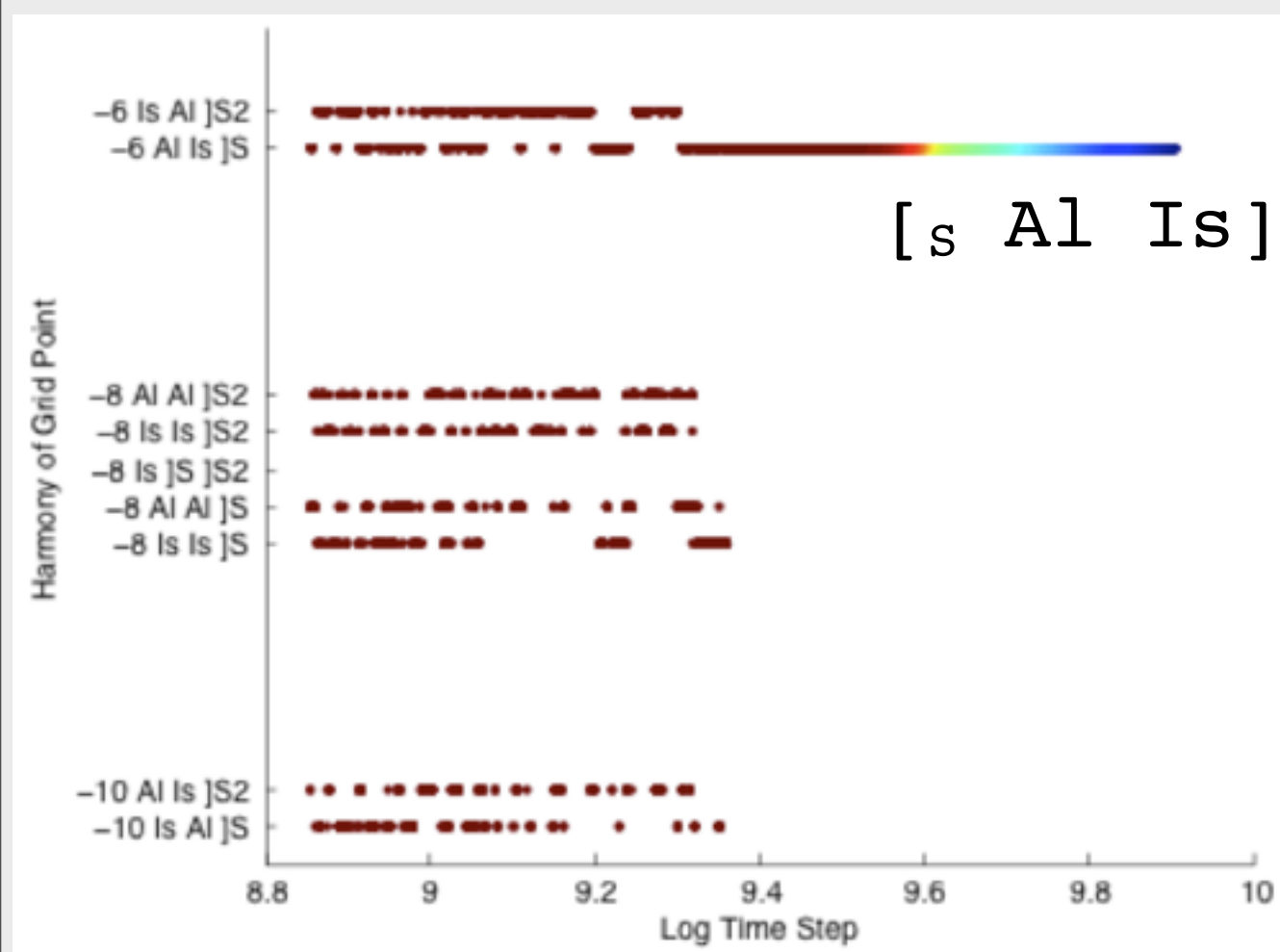
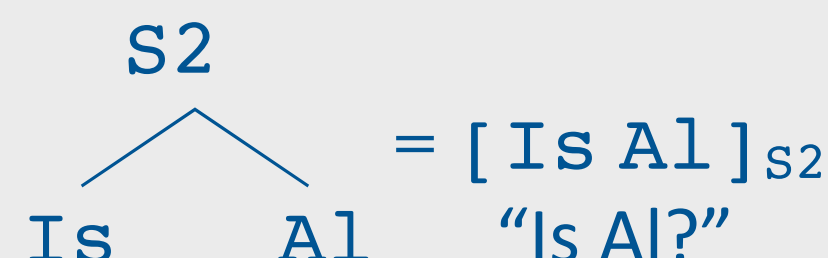
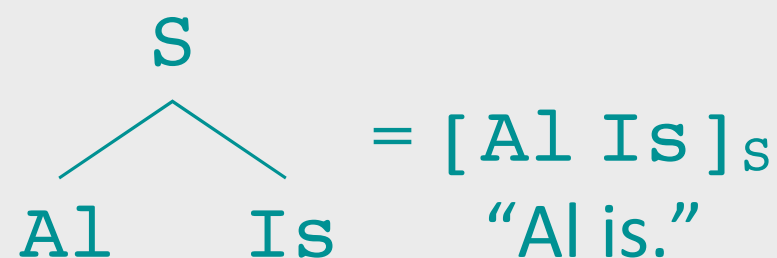
A nanogrammar \mathcal{G}

Start symbols: $\{S, S2\}$

$S \rightarrow Al\ Is$

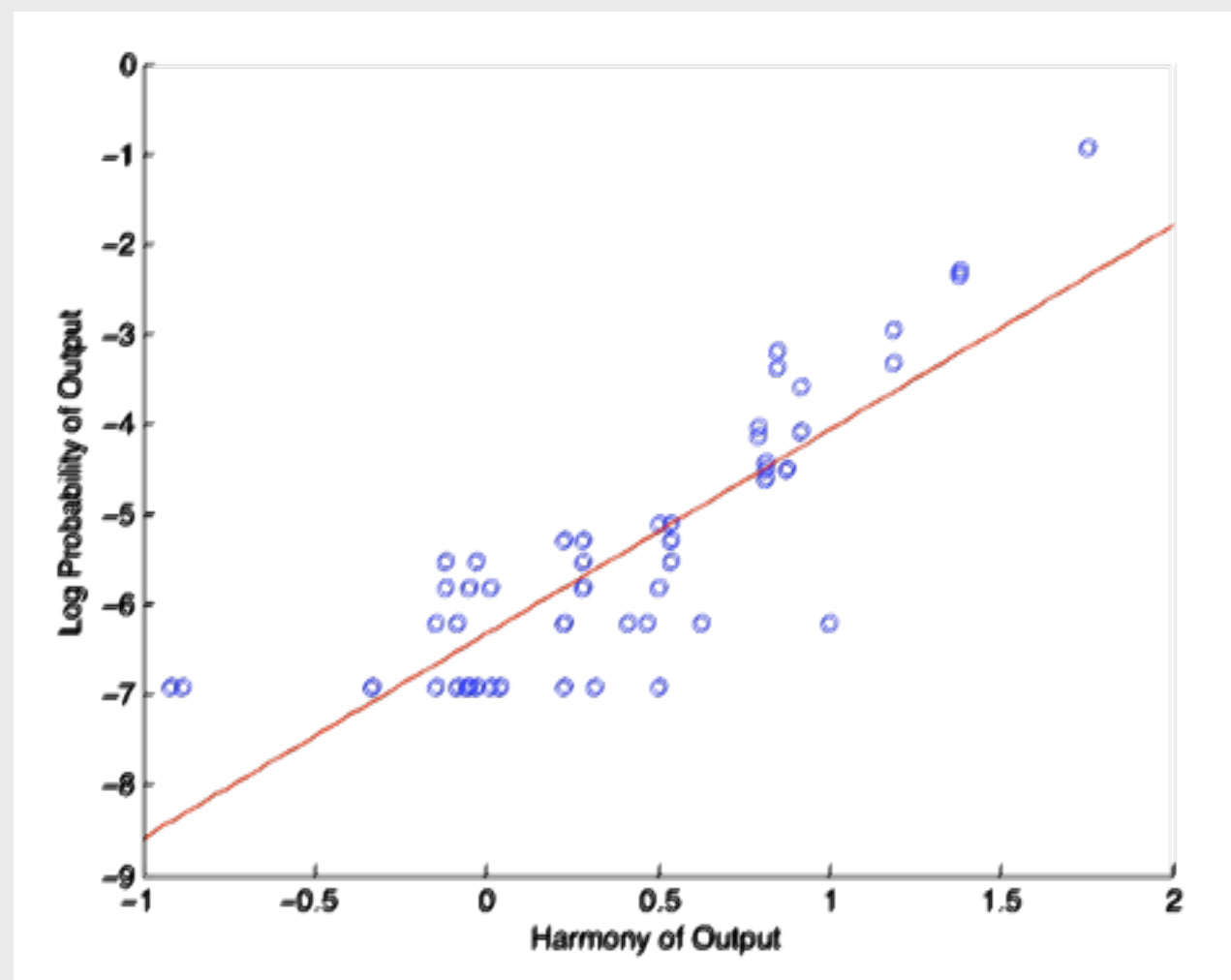
$S2 \rightarrow Is\ Al$

Its nanolanguage \mathcal{L}



Explaining error patterns with Harmony

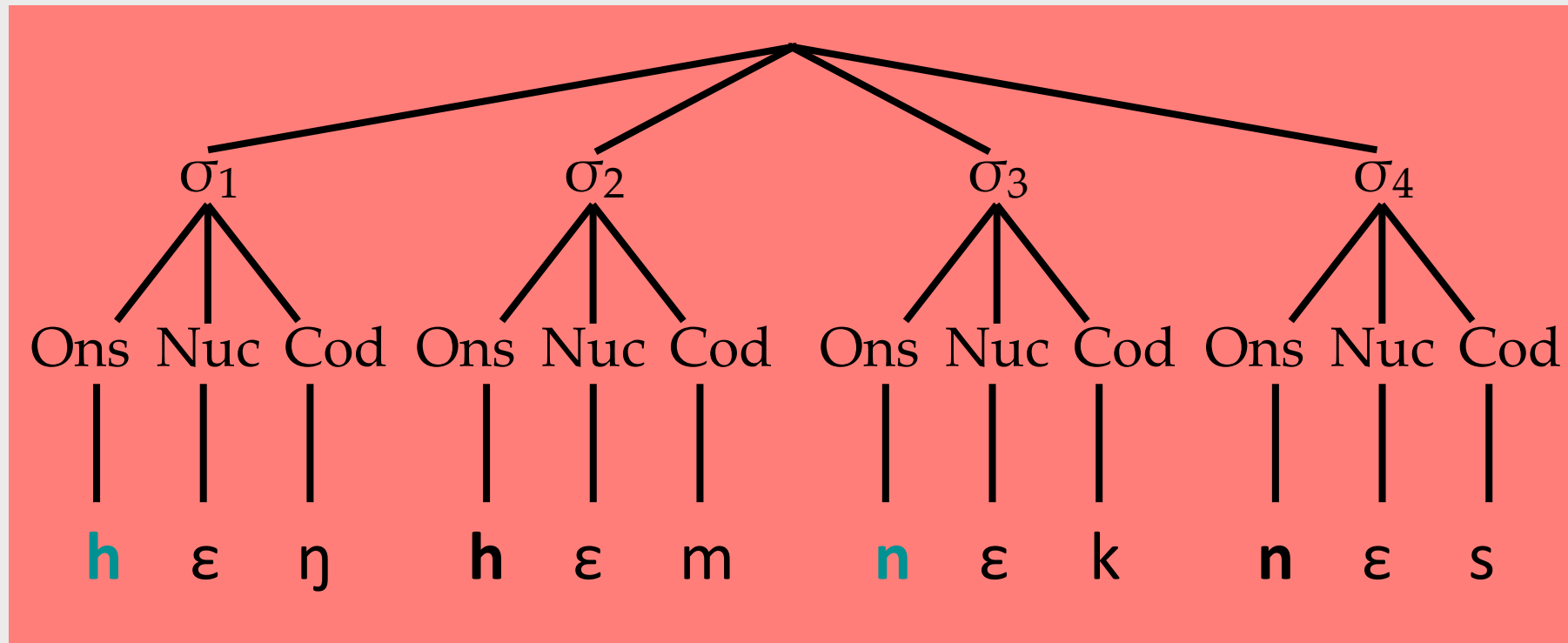
- The Harmony function is designed: we can understand it
 - ✦ H encodes the constraints of the problem domain, such that
 - ✦ the correct answer best-satisfies these constraints
- Probability of s : $p(s) \propto e^{H(s)/T}$ (T = randomness parameter)
or equivalently: $\log p(s) \propto H(s) - k$
- $/sag\ nak/ \rightarrow [?]$



Tongue-twister task
Incomplete neutralization
ITBerber syllabification

Phonological production (**g** ϵ η **h** ϵ m **f** ϵ k **n** ϵ s \rightarrow **h** ϵ η **h** ϵ m **n** ϵ k **n** ϵ s)*

- Final state of surface form component:

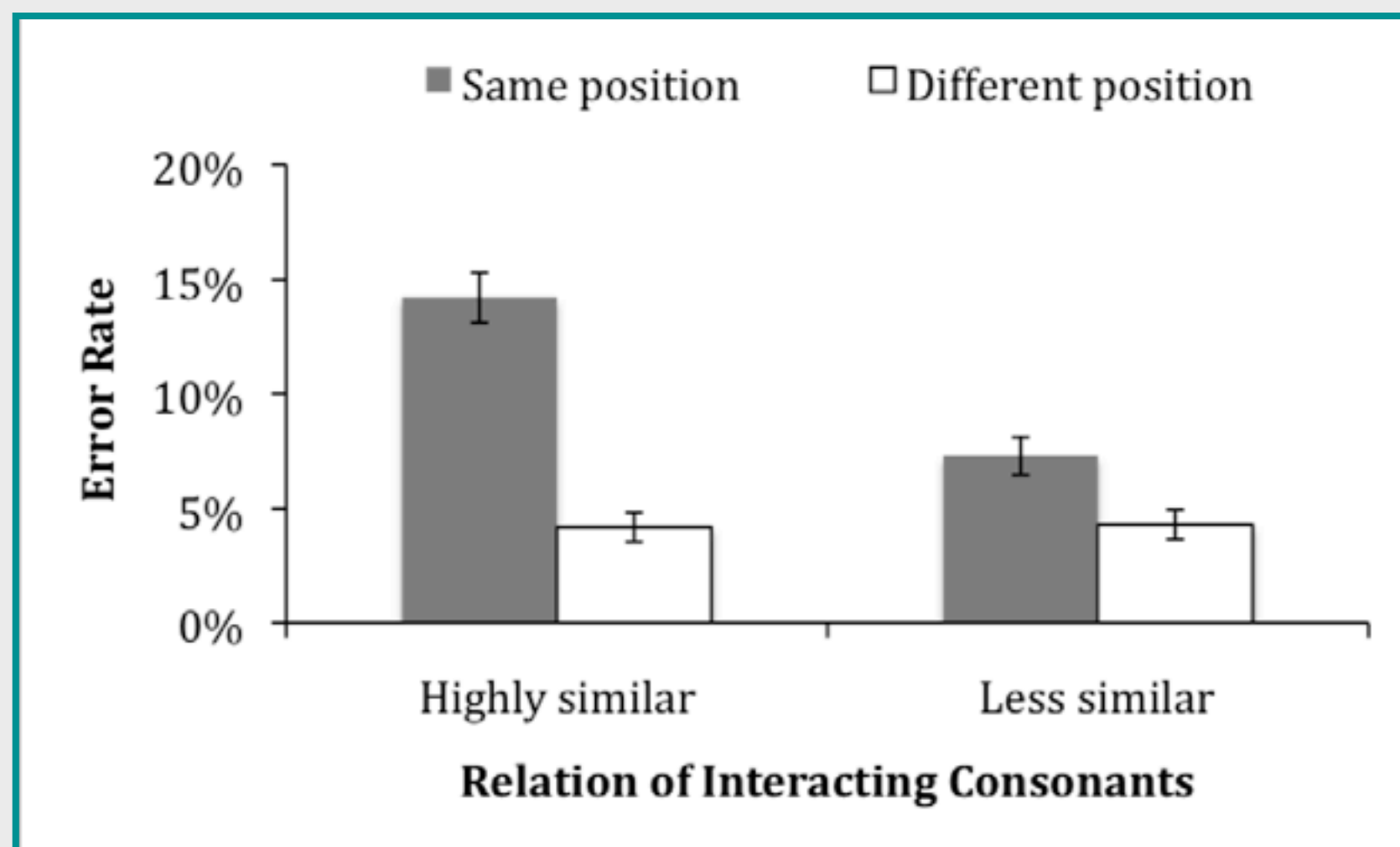


- An erroneous consonant is more likely
 - ♦ to be in a similar position (with respect to syllable structure)
 - ♦ to replace a similar consonant
 - ♦ to be in a position where it is more frequent (phonotactic probability)
 - ☞ never in a position forbidden by the English grammar (***k** ϵ **h**):
- ‘errors are well-formed’

* Dell, Reed, Adams & Meyer 2000

Phonological production (fɛŋ kɛg hɛm nɛs → fɛŋ hɛg nɛm nɛs)

- An erroneous consonant is more likely
 - ♦ to be in a similar position (with respect to syllable structure)
 - ♦ to replace a similar consonant

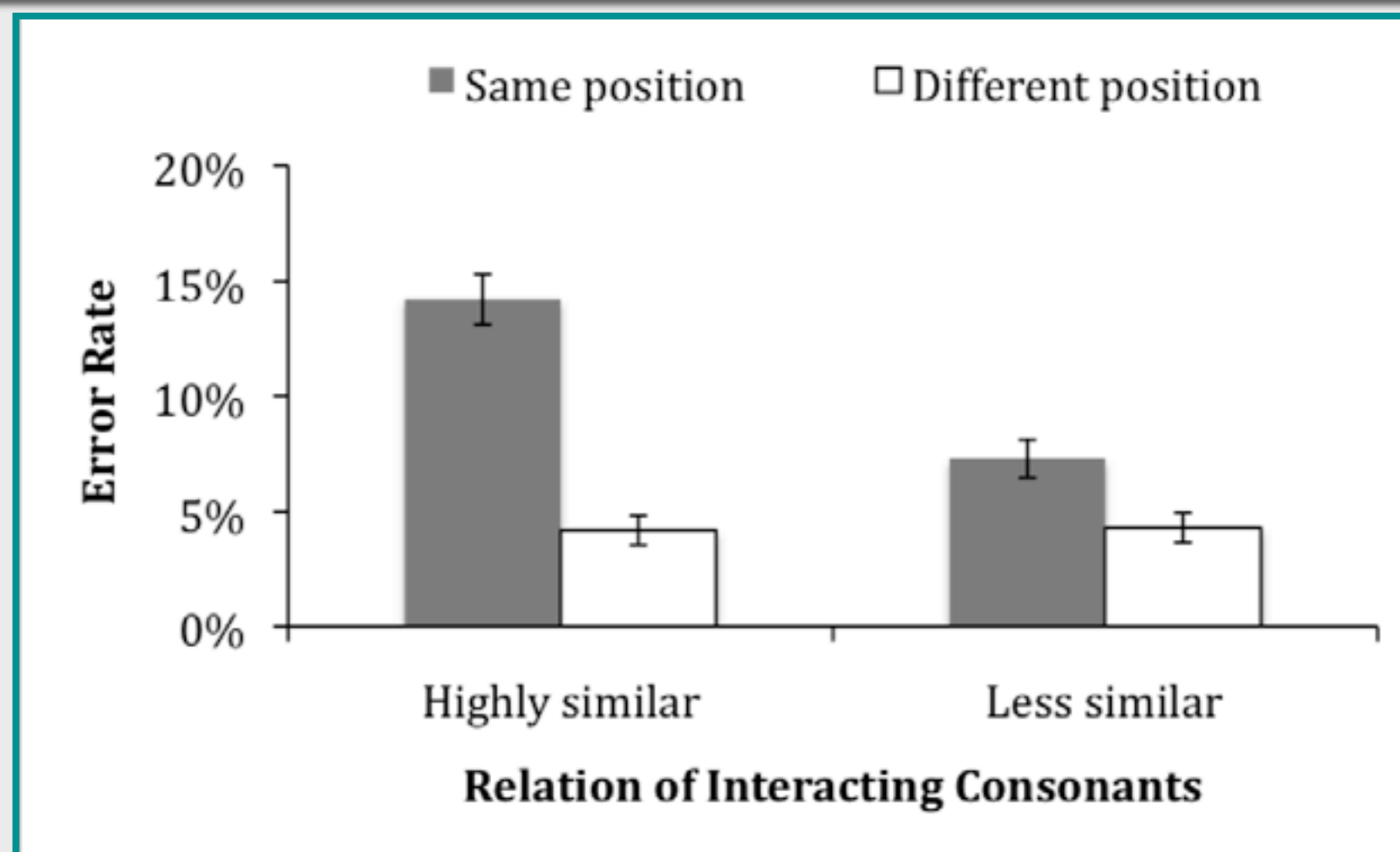


Capturing the similarity structure of roles (including recursive hierarchical structure) is a major feature of distributed tensor product representations

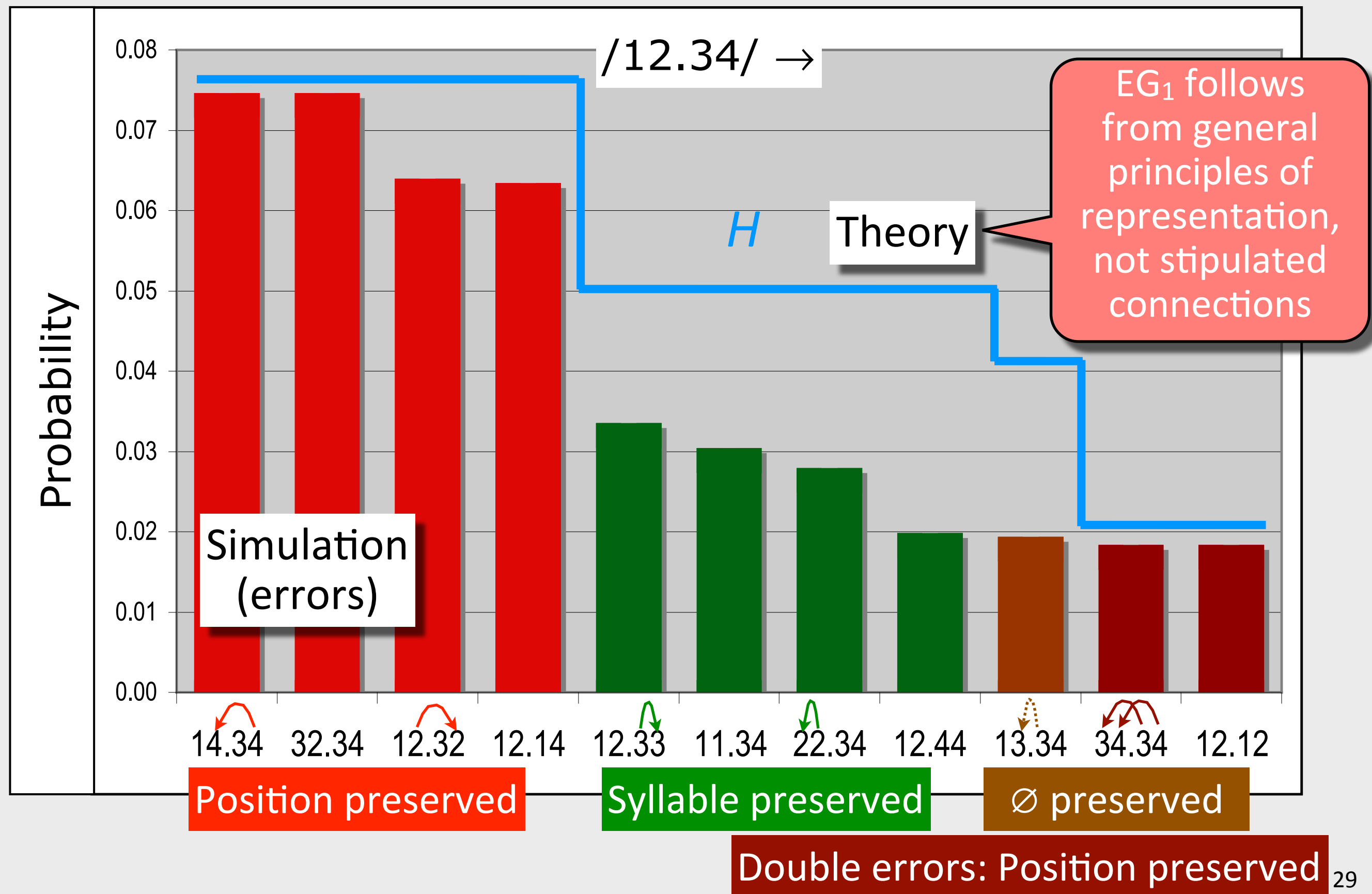
$$\mathbf{r}_{\text{Ons}/\sigma_2} = \mathbf{r}_{\text{Ons}} \otimes \mathbf{r}_{\sigma_2}$$

Preserve	Similarity	Simulations
Position	$\mathbf{r}_{\sigma_2} \cdot \mathbf{r}_{\sigma_1}$	0.4
Syllable	$\mathbf{r}_{\text{Ons}} \cdot \mathbf{r}_{\text{Cod}}$	0.1

$\text{sim}(\mathbf{r}_{\text{Ons}}, \mathbf{r}_{\text{Cod}}) < \text{sim}(\mathbf{r}_{\sigma_1}, \mathbf{r}_{\sigma_2}) \iff ?$ similarity of consonant behavior
across different positions < across different syllables]



EG₁. Errors tend to preserve position



/rad/ → [rat] 'wheel' (German)

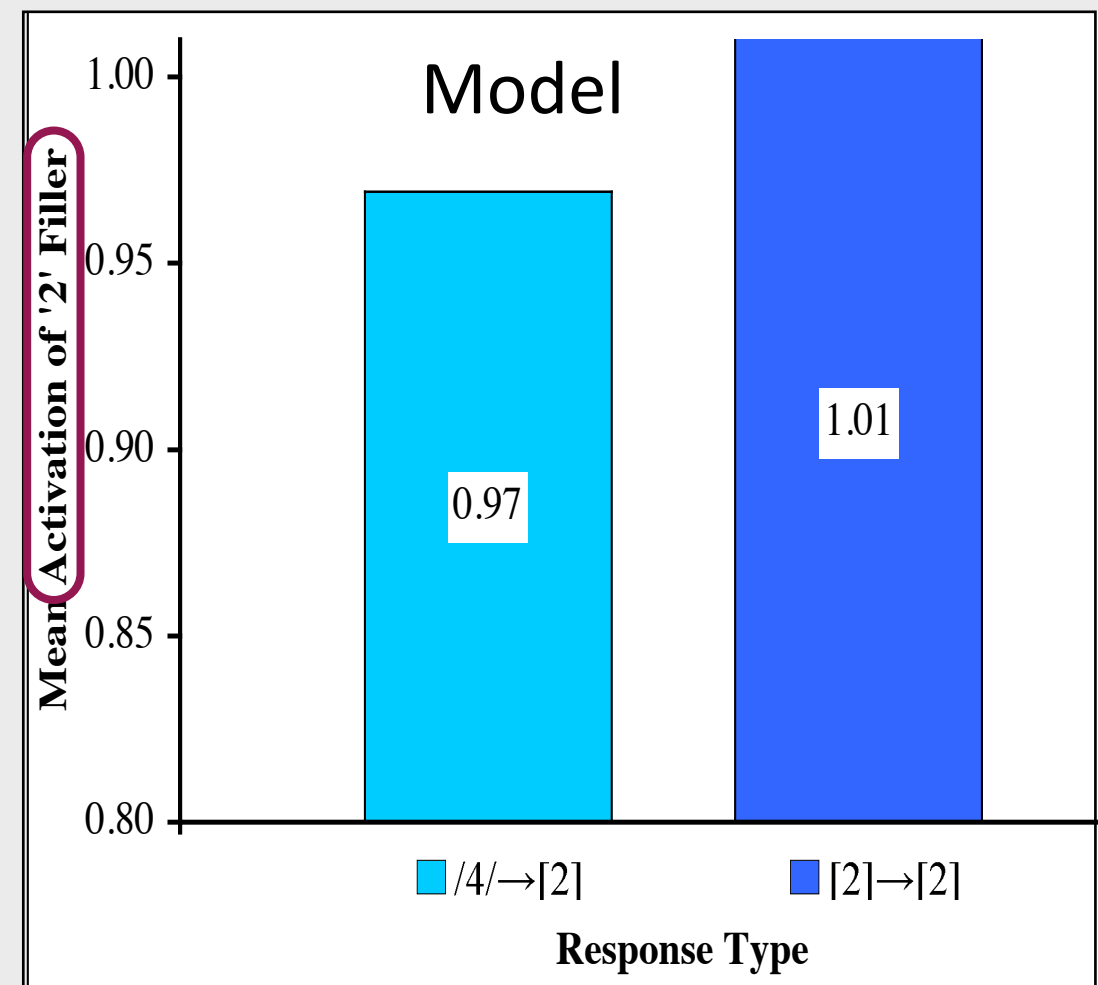
EG₇. In alternations, surface forms can show subphonemic traces of the underlying form.*

Alternations: MARK >> FAITH

E.g., **d*_{+voi}/Coda >> FAITH(voi)

Model: '4₊' vs. '2₋' ~ *d* vs. *t*

*4₊/Coda (1.5) >> FAITH (1.0)



Surface forms can show subphonemic 'traces' of the underlying form.

Review: Warner, Jongman, Sereno & Kemps, 2004 (*Journal of Phonetics*)

* Factors other than underlying form can also induce similar effects.

Model: derived *directly* from the HG
version of the OT analysis of P&S93

