# Using Artificial-Intelligence-Driven Deep Neural Networks to Uncover Principles of Brain Representation and Organization

2017.04.27

## Daniel Yamins

Stanford Neurosciences Institute
Stanford Artificial Intelligence Laboratory
Departments of Psychology and Computer Science
Stanford University

Understanding complex, noisy data streams is a critical part of cognition.

Understanding complex, noisy data streams is a critical part of cognition.



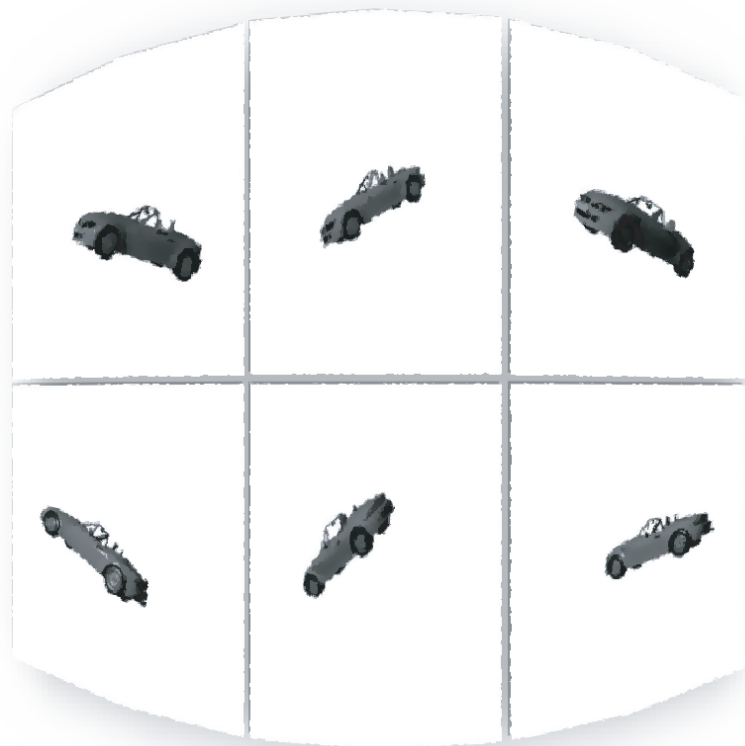"Mercedes behind Lamborghini, on a field in front of mountains."

# Variation Makes Object Recognition Challenging



View: position, size, pose, illumination
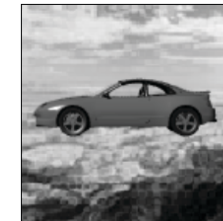
Background variation

Distortion & Noise

Geometric variation

Beetle　　BMW Z3　　Clio
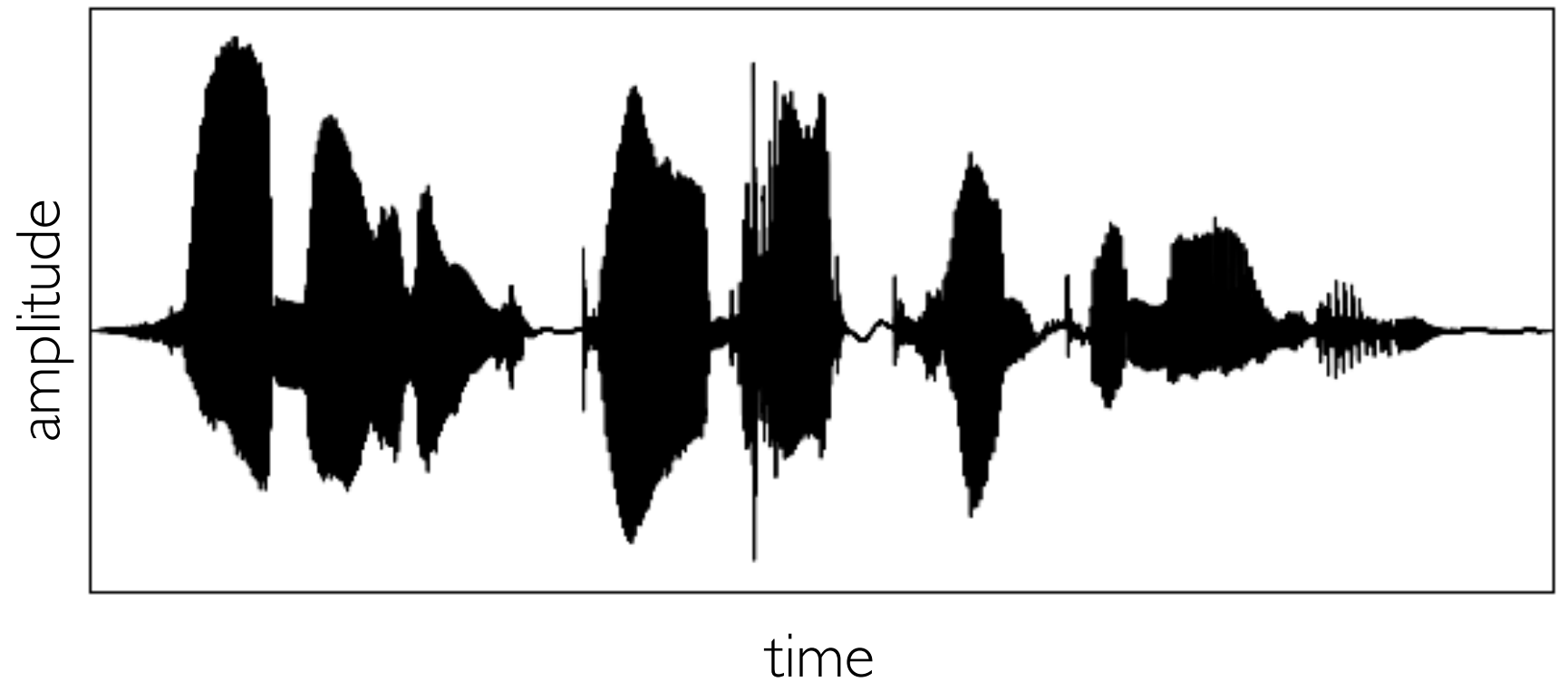
Celica　　　　　　　　Alfa

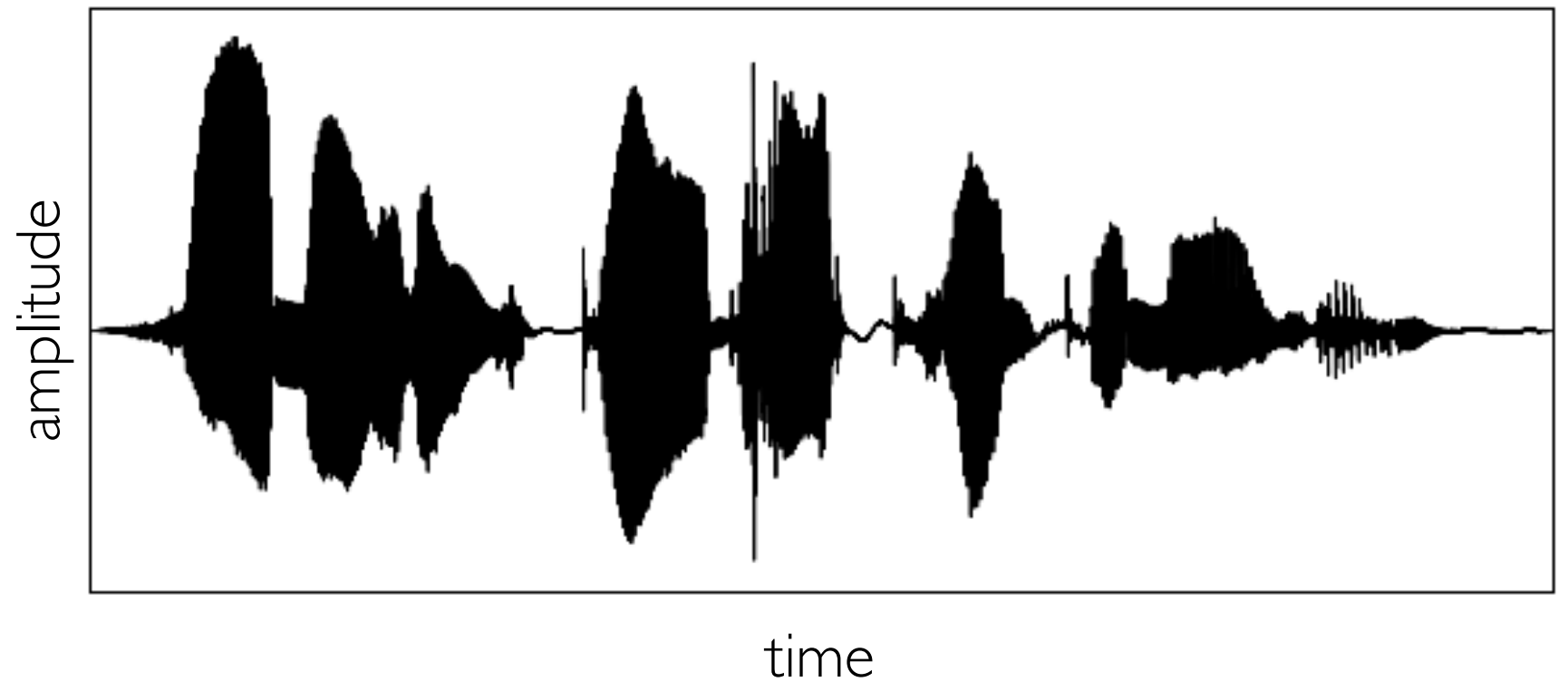car identities

VW Bora　BMW 325　Astra

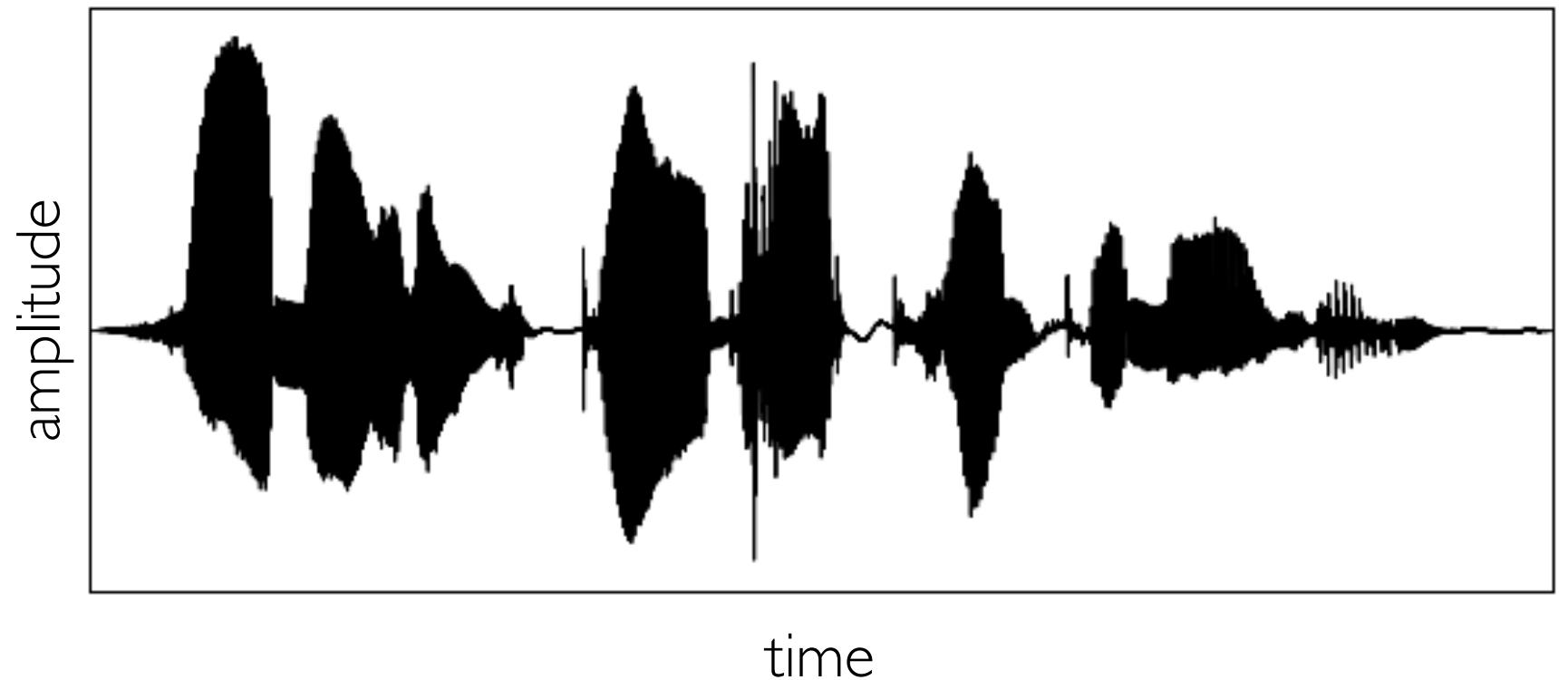Understanding complex, noisy data streams is a critical part of cognition.



amplitude

time

Understanding complex, noisy data streams is critical part of cognition.



"Hannah is good at compromising."

Understanding complex, noisy data streams is critical part of cognition.



"Hannah is good at compromising."

variation sources:  speaker identity

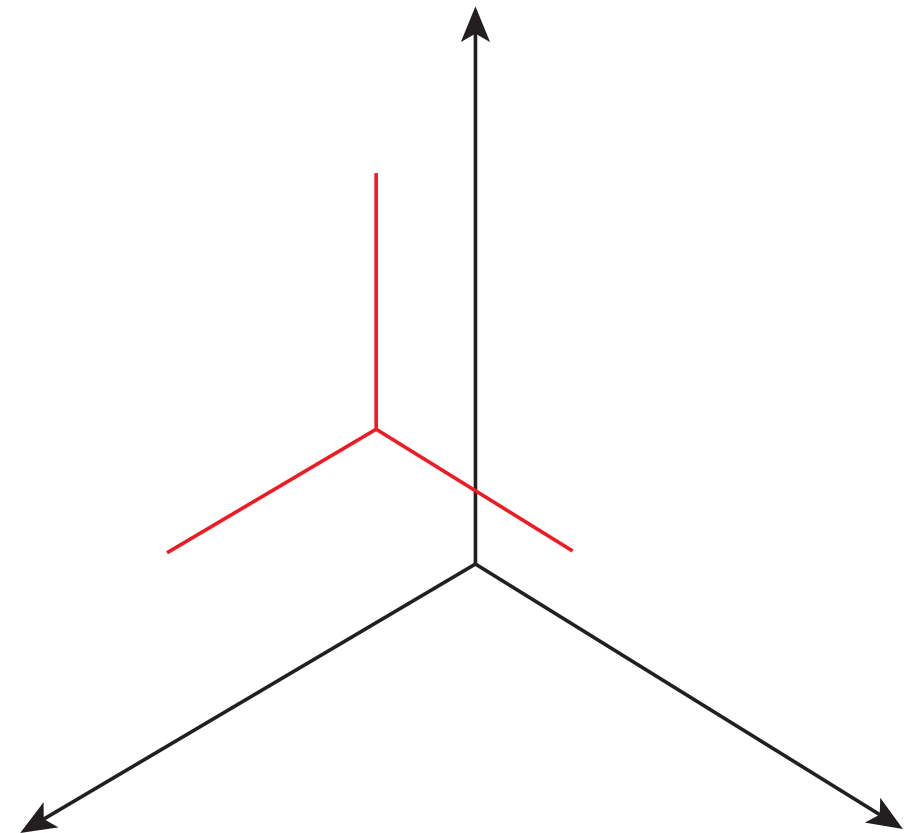background noise

reverberation

…

Axes of natural variation of natural
**"physics"** representation of world

e.g.

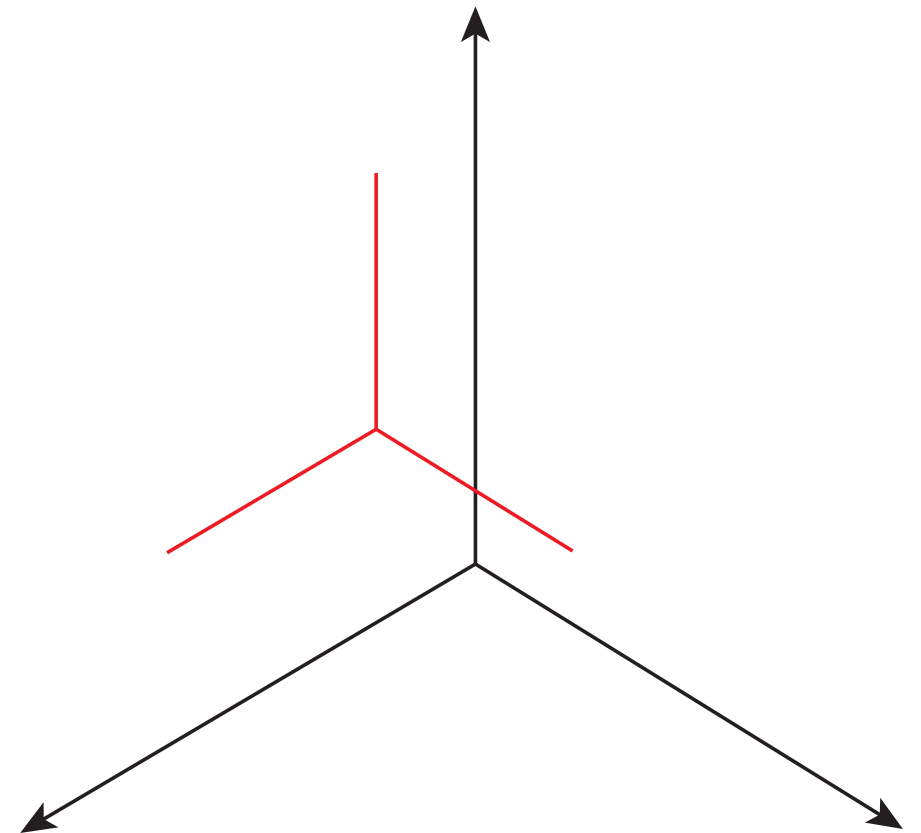retinal photoreceptor voltage
or hair-cell point amplitudes

Axes of natural variation for
natural **behavioral** events
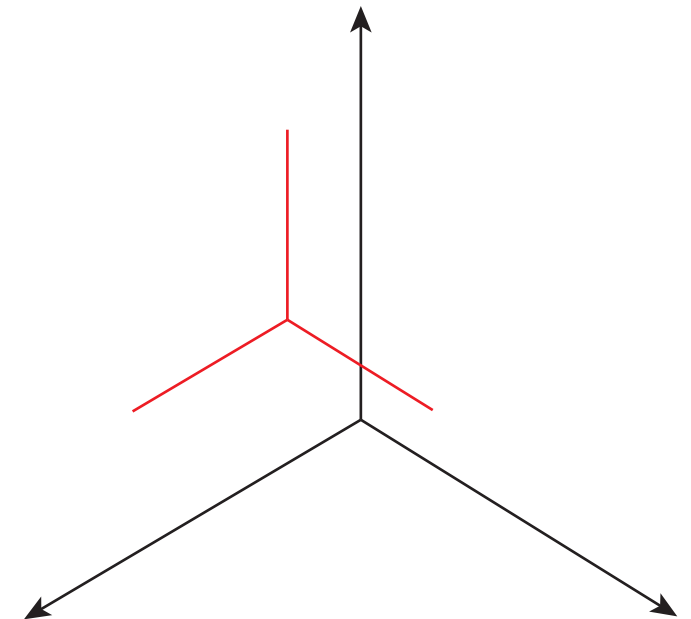
e.g.

deforming face moving in complex-
lighted environment
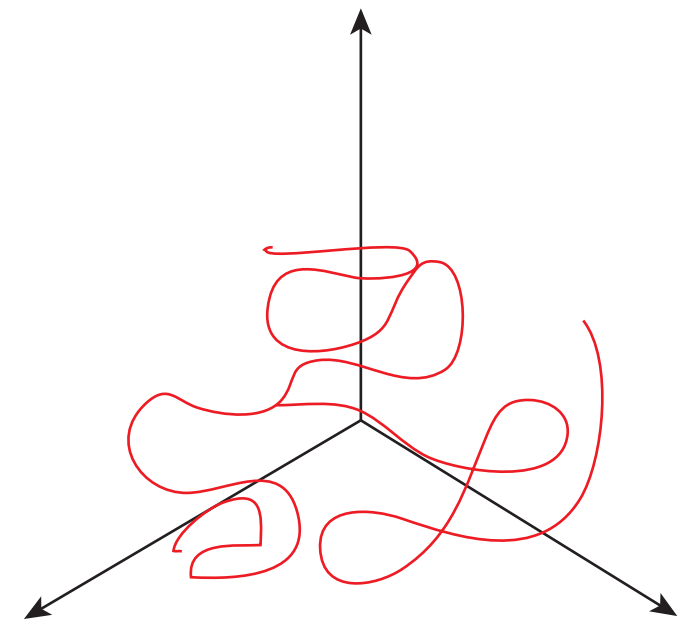
Axes of natural variation for
natural **behavioral** events
(e.g. deforming face moving in
complex-lighted environment)

*are misaligned with*

Axes of natural variation of natural
**"physics"** representation of world
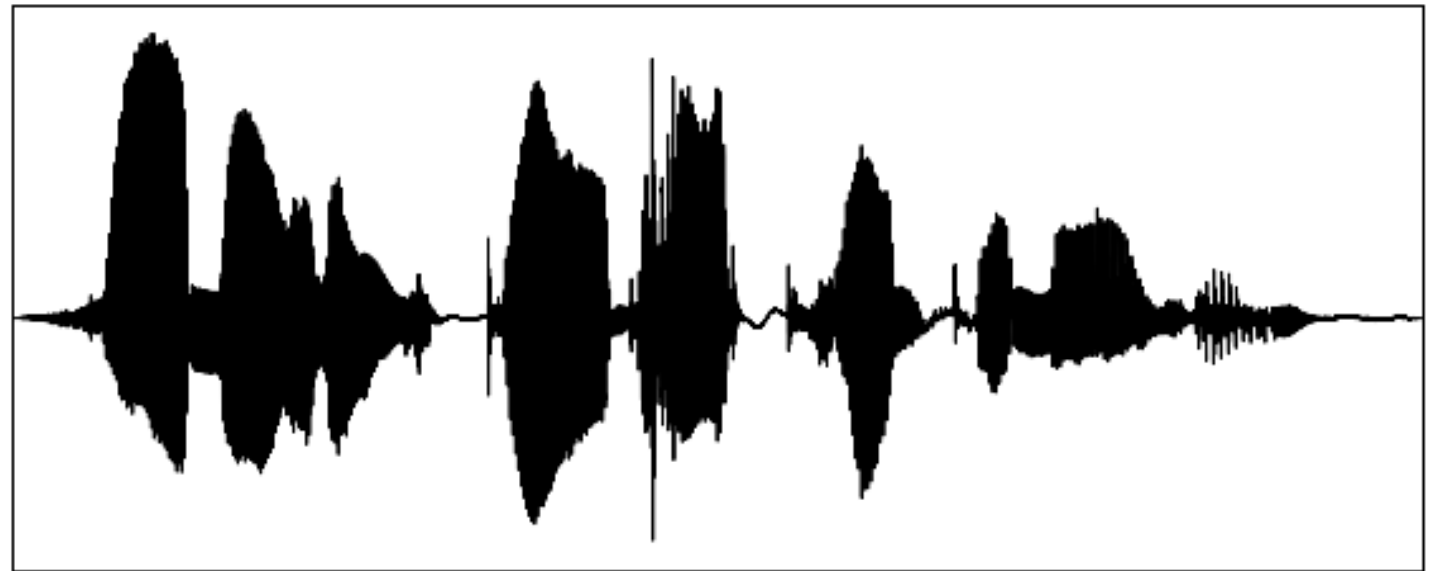e.g. retinal photoreceptor voltage
or hair-cell point amplitudes

visual cortex

auditory cortex

"Mercedes behind Lamborghini, on a field in front of mountains."

"Hannah is good at compromising"

"Mercedes behind Lamborghini, on a field in front of mountains."

"Hannah is good at compromising"

visual cortex

auditory cortex

"Mercedes behind Lamborghini, on a field in front of mountains."
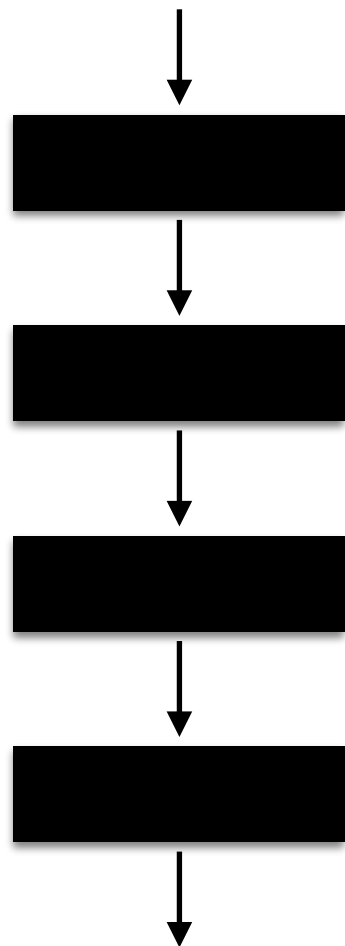
"Hannah is good at compromising"

# Ventral Stream = Connected series of brain areas
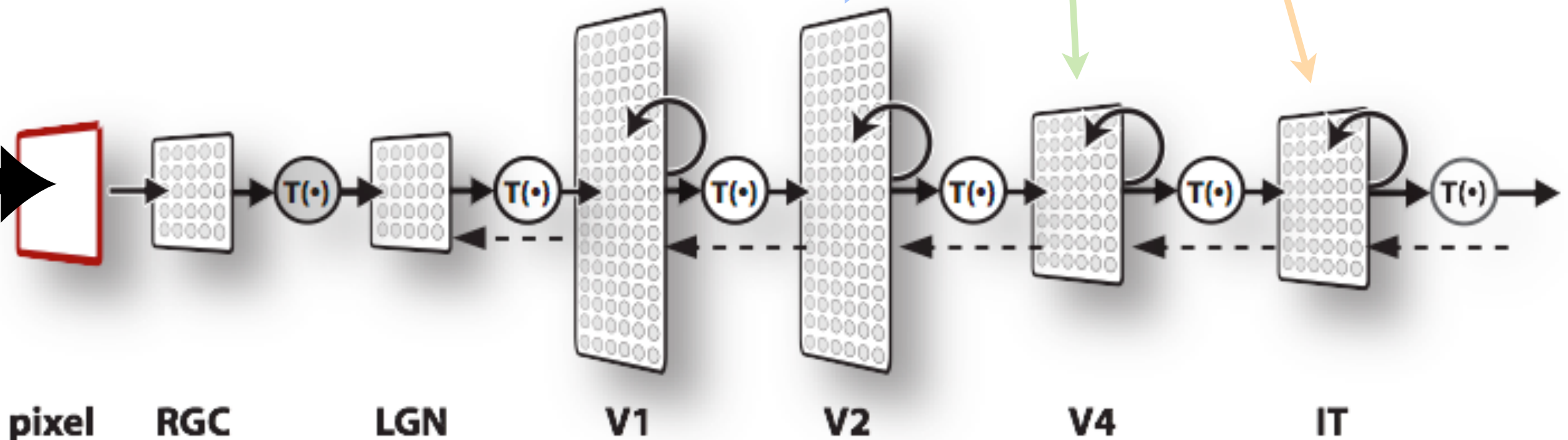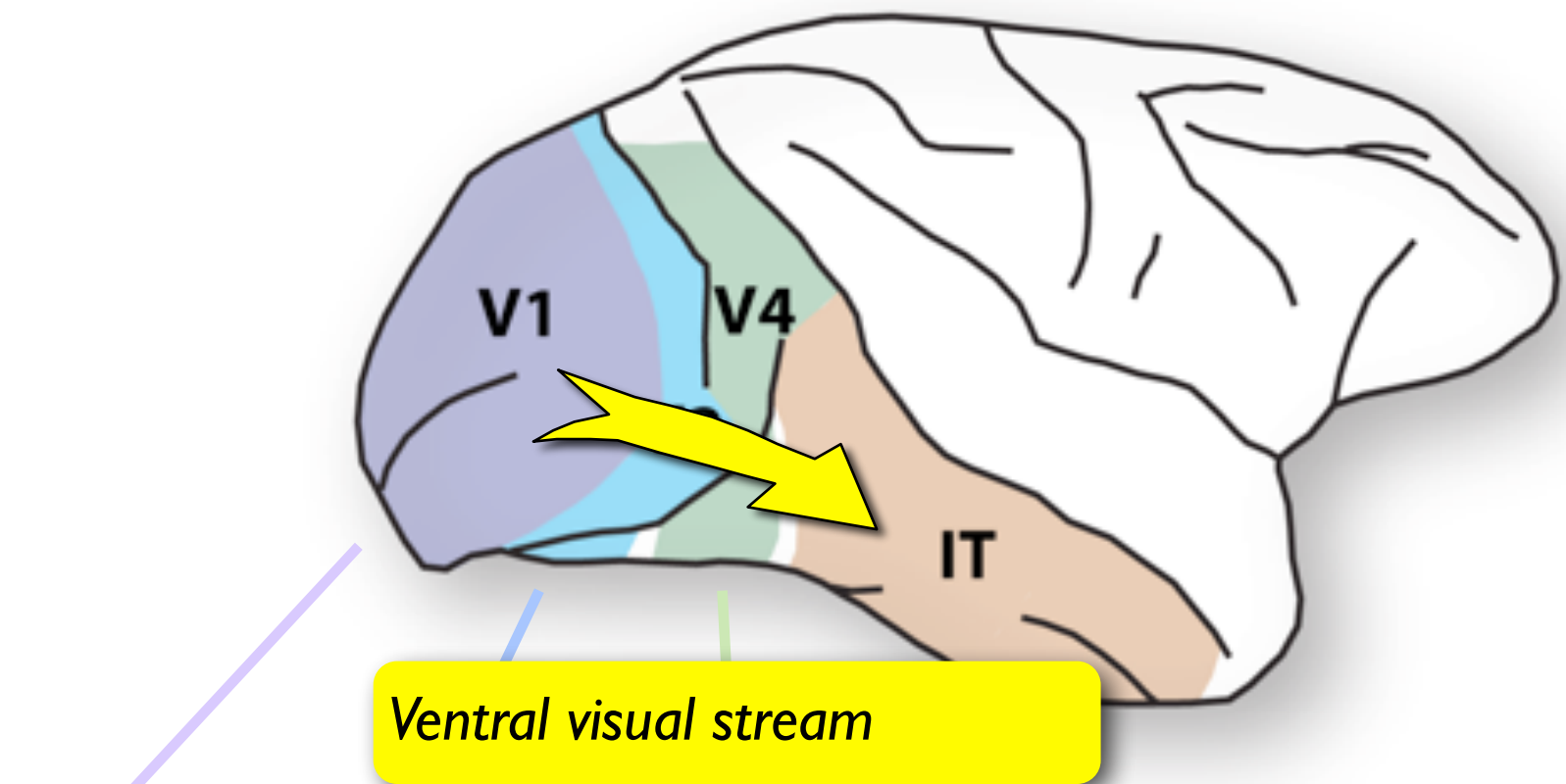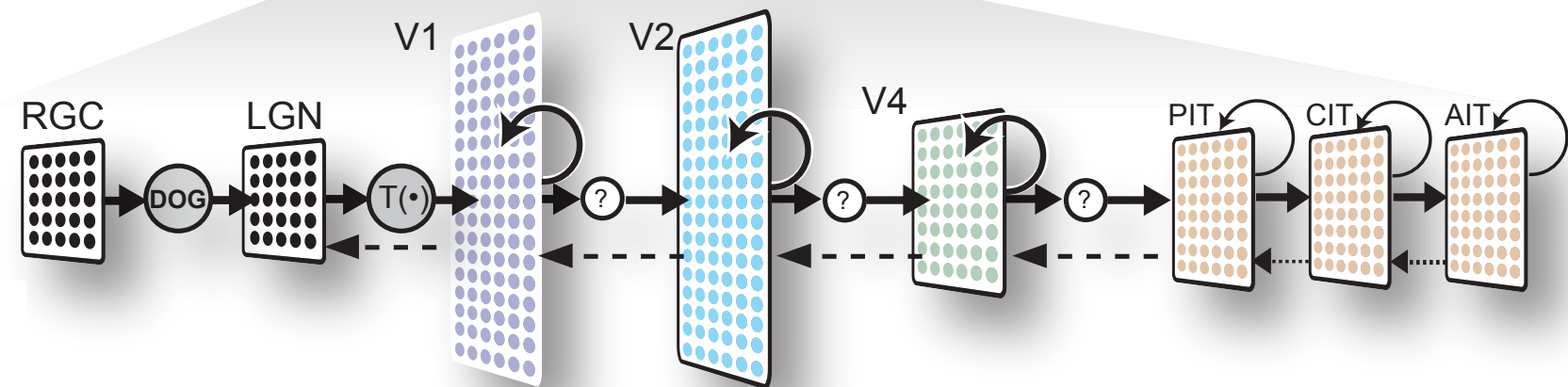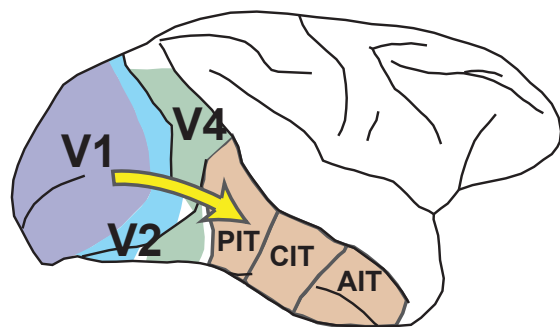
Kaas (2003), Van Essen (2003), Valois and Morgan (1974) Gross (1973), Mishkin and Ungerleider (1983), Holmes and Gross (1984) Horel *et al.* (1987); Freiwald and Tsao (2010), Pitcher, *et al.* (2009) Yaginuma (1982), Holmes (1984), Weiskrantz (1984), Schiller (1995)  Afraz (2006), Verhoef (2012)  Rust (2010), Freiwald (2010), Lehky (2007) Majaj (2015)
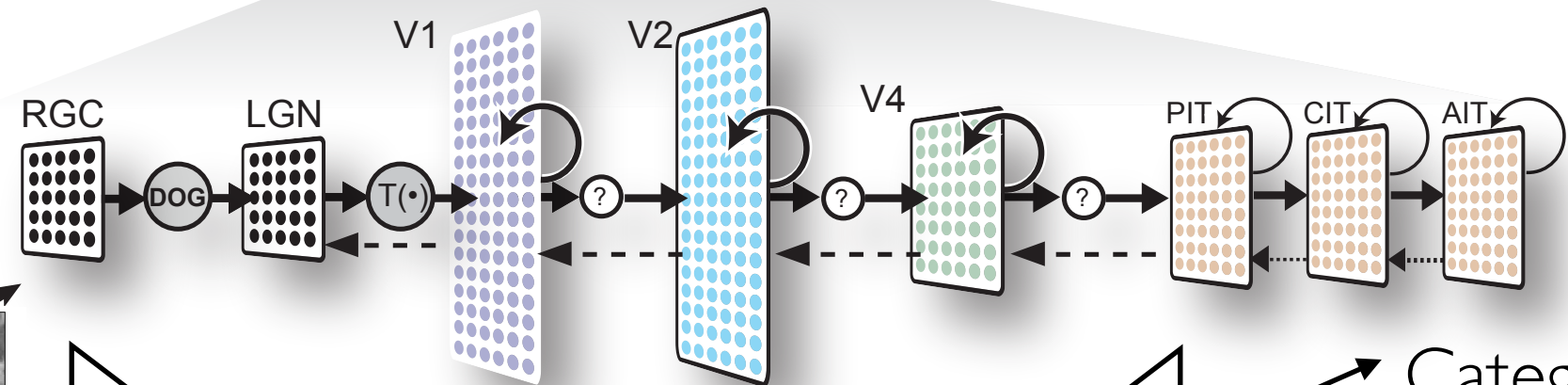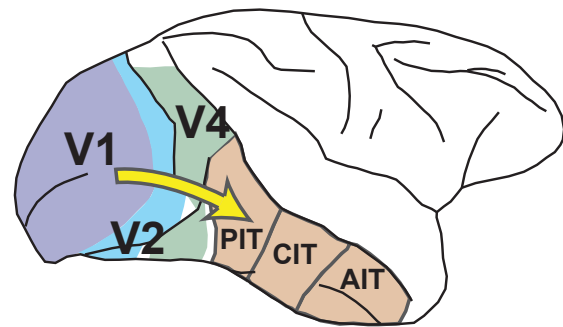
**pixel**  **RGC**  **LGN**  **V1**  **V2**  **V4**  **IT**

Stimulus → *representation* → Neurons → *read-out* → Behavior

Stimulus $\xrightarrow{\textit{representation}}$ Neurons $\xrightarrow{\textit{read-out}}$ Behavior
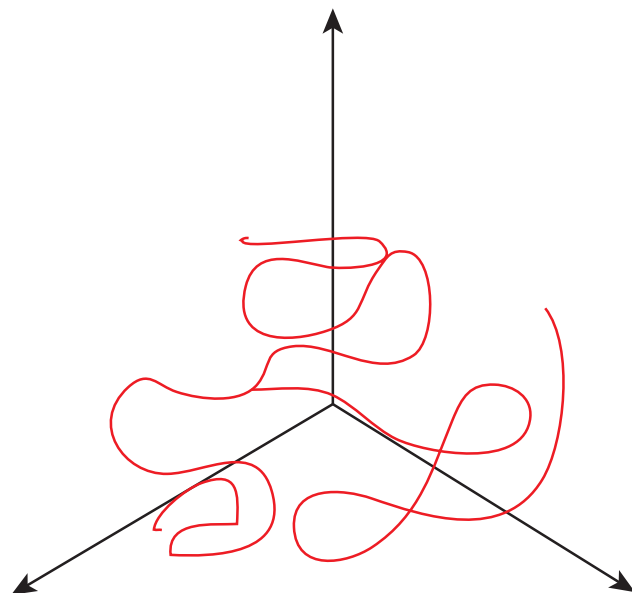
V1   V2   V4

RGC   LGN   PIT   CIT   AIT

visual representation

Category
Location
Size
Pose
Depth relationships

Multi-array electrophysiology in macaque V4 and IT.



V4

IT

10mm

☐ = Array

# Multi-array Electrophysiology Experiment

5760 images

64 objects

8 categories

uncorrelated photo backgrounds

Low variation

... *640 images*

Medium variation

... *2560 images*

High variation

... *2560 images*

Animals   Boats   Cars   Chairs   Faces   Fruits   Planes   Tables

Pose, position, scale, and background variation

# Multi-array Electrophysiology Experiment

About 300 total sites

= Array

V4

IT

10mm

*Output = Binned spike counts 70ms-170ms post stimulus presentation averaged over 25-50 reps of each image.*

Img 1   Img 2   Img **5760**

Neuron 1
Neuron 2
Neuron 3
⋮
Neuron **296**

...

...

# Neural-Behavior Decoding



Img 1    Img 2    Img **5760**

Animal or not?

linear combination of units

Neuron 1
Neuron 2
Neuron 3
⋮
Neuron **296**

different linear combination

Car or not?

Chair or not?

Face or not?

V4 loses out at higher variation:



Basic
categorization

# Decoding Behaviorally Output from Neural Populations

V4 loses out at higher variation:

… but humans are much less affected.

Basic categorization



*Yamins\* and Hong\* et. al. **PNAS** (2014)*

# IT Neurons Track Human Performance

V4 loses out at higher variation:

… but humans are much less affected.

… as is the IT neural population.

Basic categorization



*Yamins\* and Hong\* et. al.* **PNAS** *(2014)*

# IT Neurons Track Human Performance

V4 loses out at higher variation:

… but humans are much less affected.

… as is the IT neural population.

Basic categorization

Yamins* and Hong* et. al. **PNAS** (2014)

At <u>high variation levels</u>, IT much better than V4 and existing models.

# IT Neurons Track Human Performance

IT matches human error patterns as well as raw performance.



IT Population

V4 Population

Human Performance

Neural Decode Performance

Human Dprime

Neural D-prime (128 units)

● Low-Variation Face subordinate tasks.

**Human**                **Rhesus monkey**

Camel
Dog
Rhino
Elephant
Wrench
Knife
Hanger
Fork
Guitar
Pen
Tank
Truck
Bird
Hammer
Gun
Table
Calculator
Spider
Leg
Zebra
House
Bear

"camel" confused with "dog"

"tank" confused with "truck"

**Upshot: human and non-human primate basic level core object percepti (sp. identification) are indistinguishable**

**Does not depend on reporting effector (touch vs. eye movement)**

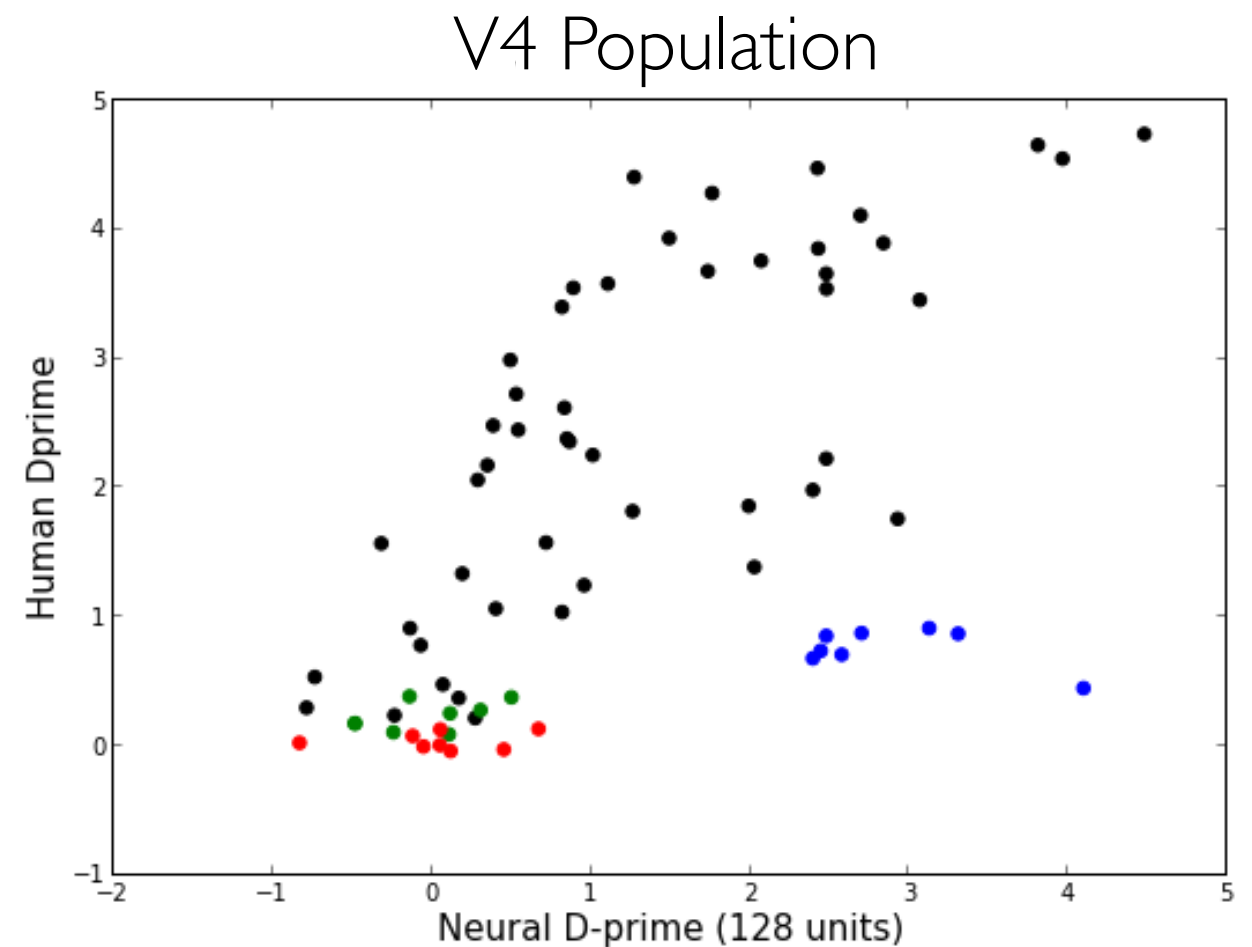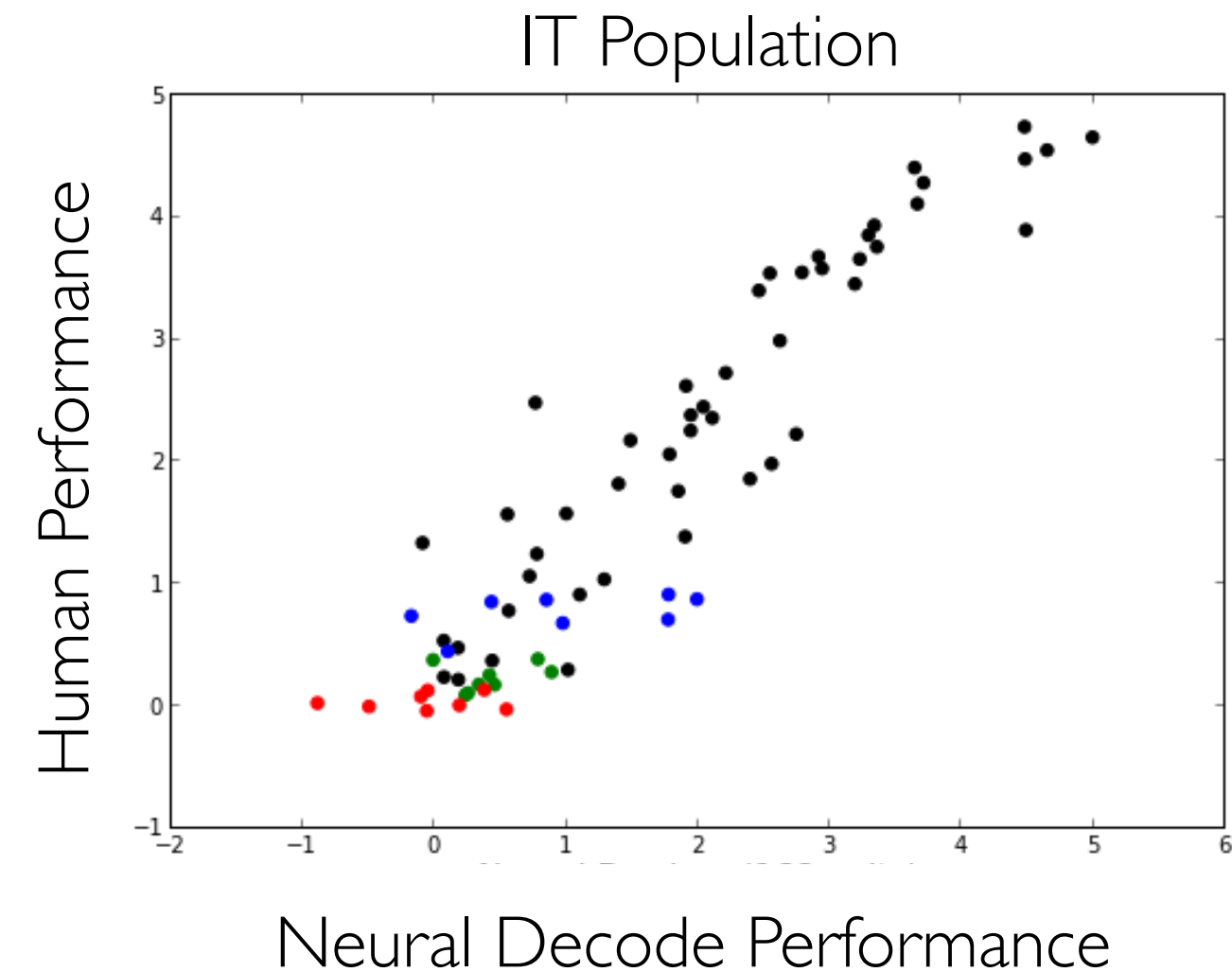Comparison of Object Recognition Behavior in Human and Monkey
R. Rajalingham, K Schmidt, J.J. DiCarlo, **Vision Sciences Society** (2014)
R. Rajalingham, K Schmidt, J.J. DiCarlo, **J. Neuroscience** (2015)

Adapted from Motter and Mountcastle 1981

**GOAL: Predictive model of single-neuron responses throughout the ventral stream to arbitrary image stimuli.**

**GOAL: Predictive model of single-neuron responses throughout the ventral stream to arbitrary image stimuli.**



1. image-computable          2. Predictive          3. Mappable

**GOAL:  Predictive model of single-neuron responses throughout the ventral stream to arbitrary image stimuli.**



pixel    RGC      LGN      V1      V2      V4      IT
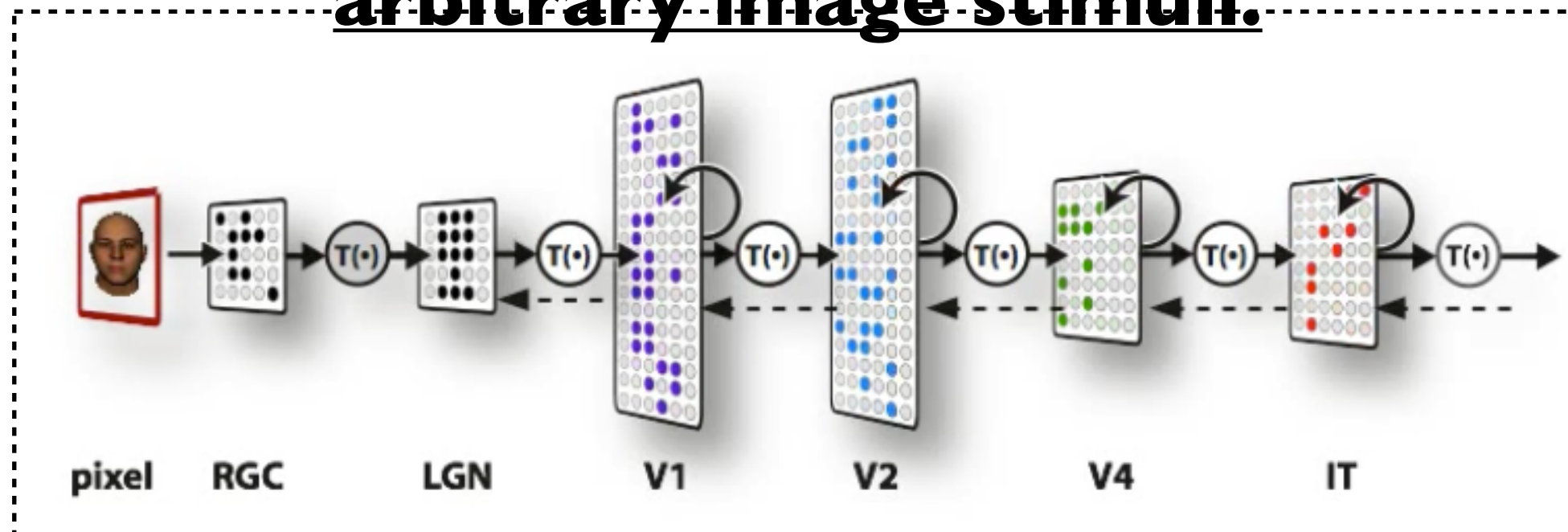
1. image-computable      2. Predictive      3. Mappable

→   Convolutional Neural Networks (CNNs)

Fukushima, 1980; Lecun, 1995

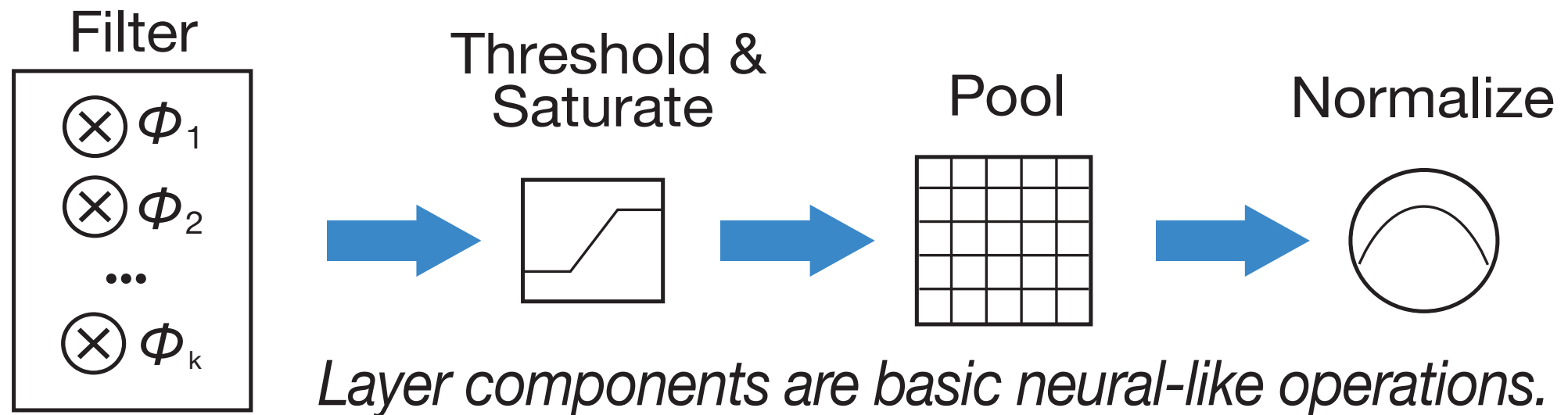Kunihiko Fukushima!

*Tokyo, November 2015*

Kunihiko Fukushima!

Developed neocognitron while Japan Broadcasting Corporation (NHK) … office directly next door to Keisuke Toyama and Keiji Tanaka
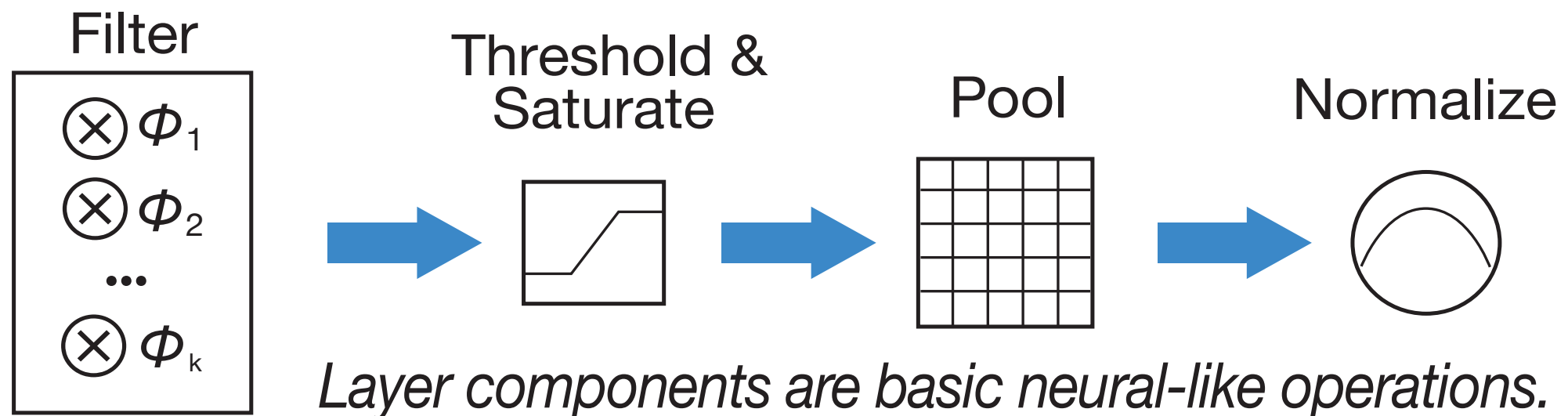
*Tokyo, November 2015*

▶ Individual layers of neurally-plausible **basic operations**



*Layer components are basic neural-like operations.*

▶ Individual layers of neurally-plausible **basic operations**

| Filter | Threshold & Saturate | Pool | Normalize |
|--------|----------------------|------|-----------|
| $\otimes \Phi_1$ $\otimes \Phi_2$ ... $\otimes \Phi_k$ | | | |

*Layer components are basic neural-like operations.*

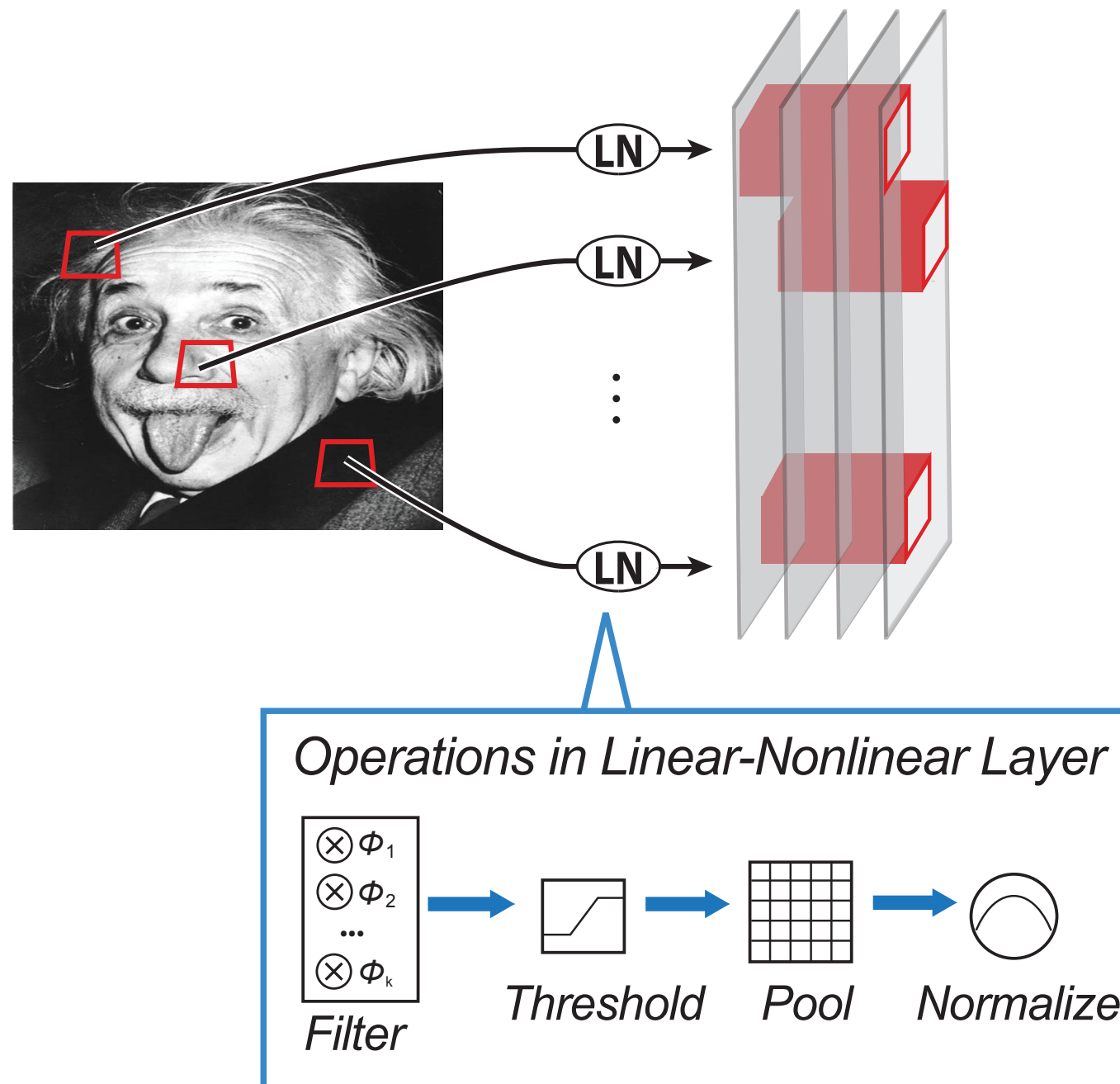|  | Filter | Threshold & Saturate | Pool | Normalize |
|--|--------|----------------------|------|-----------|
| **neuro:** | synaptic weights patterns | single-unit activations | complex cells | competitive inhibition |
| **data:** | untangling through dimension expansion | "AND" operation by limiting dynamic range | adding robustness by dimension reduction | put results back into standard range |

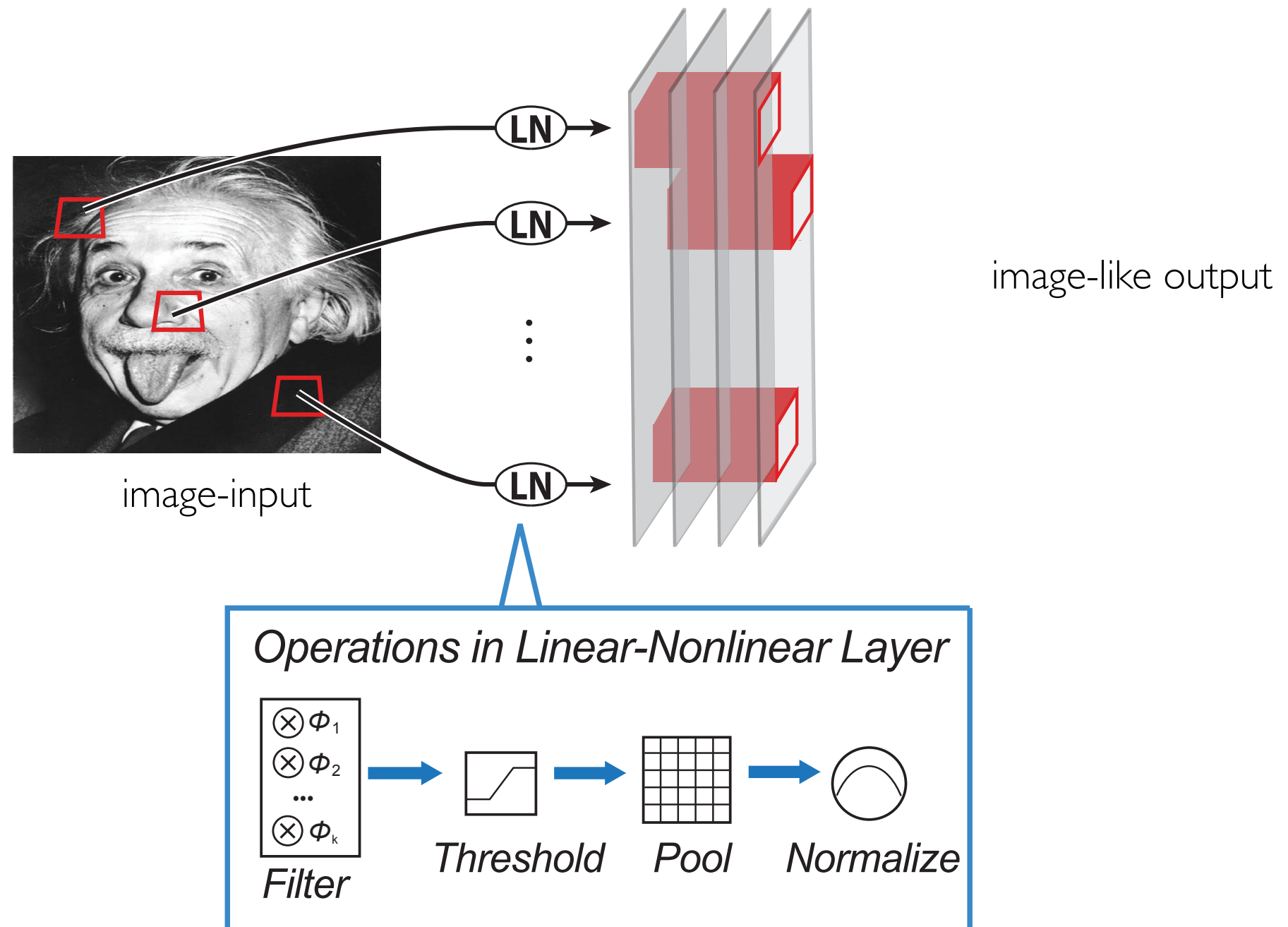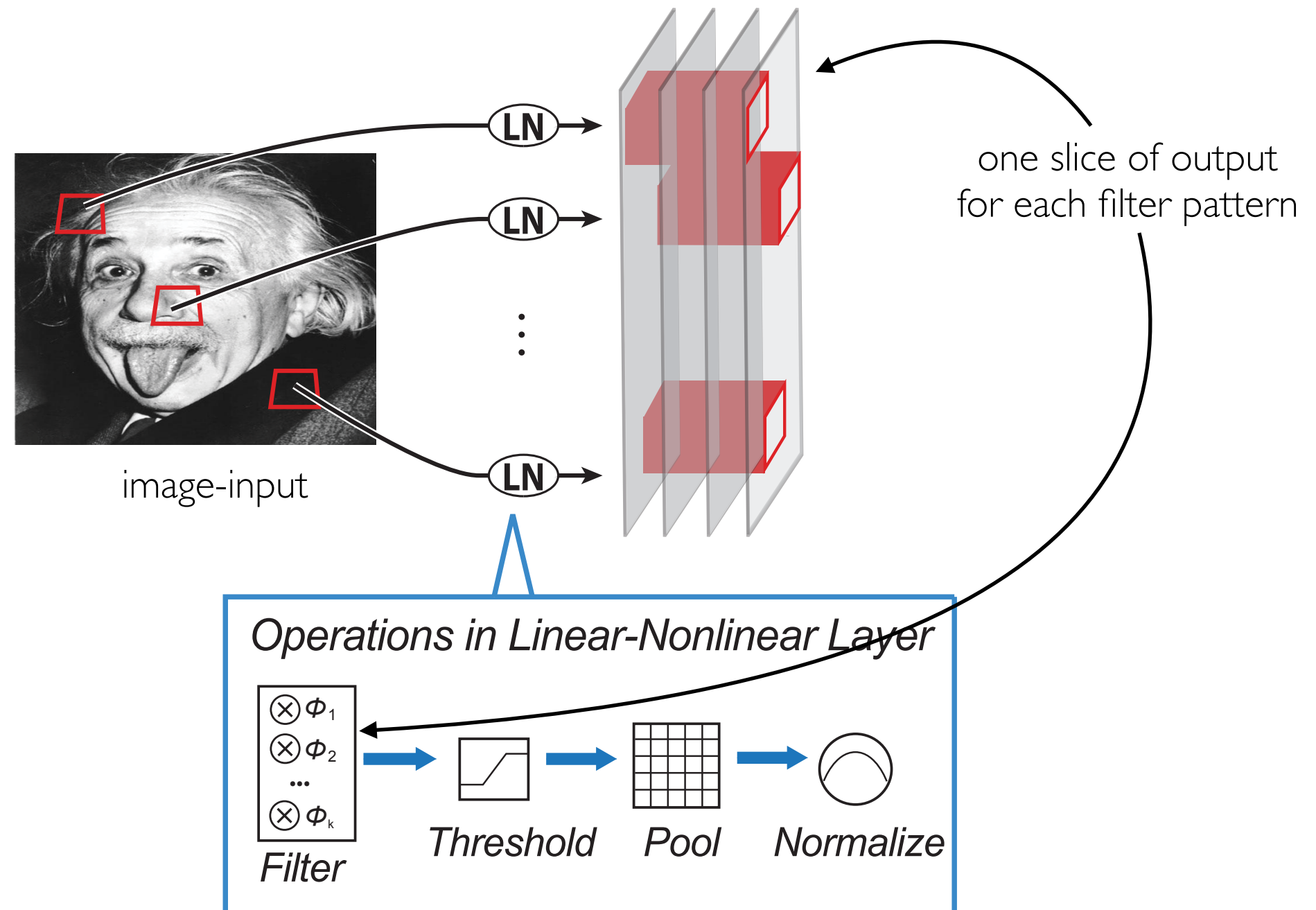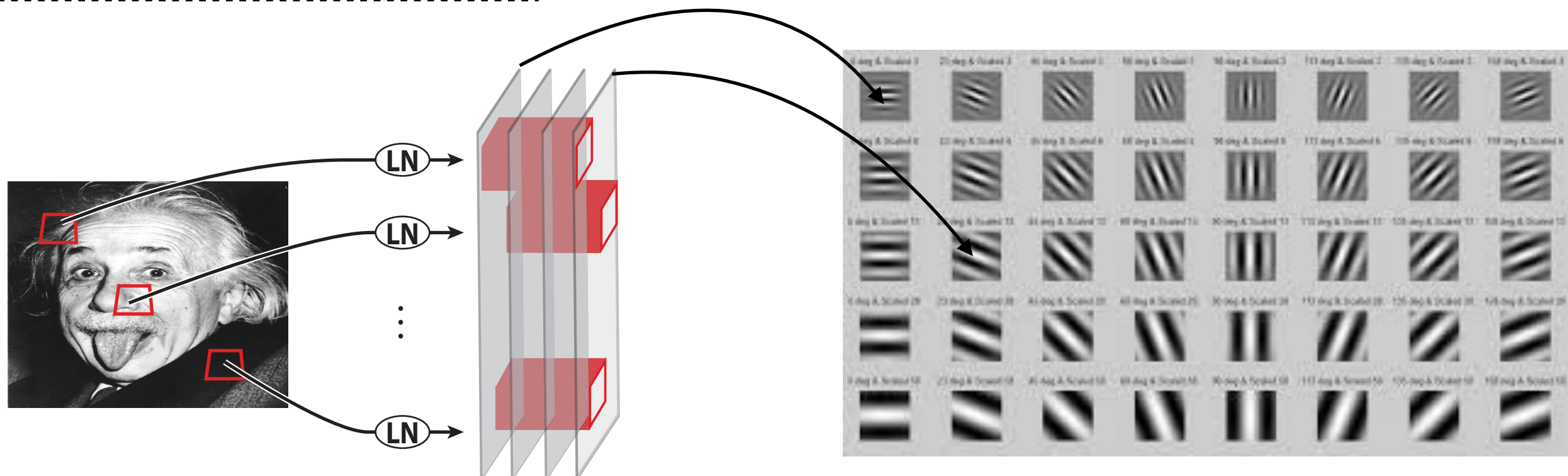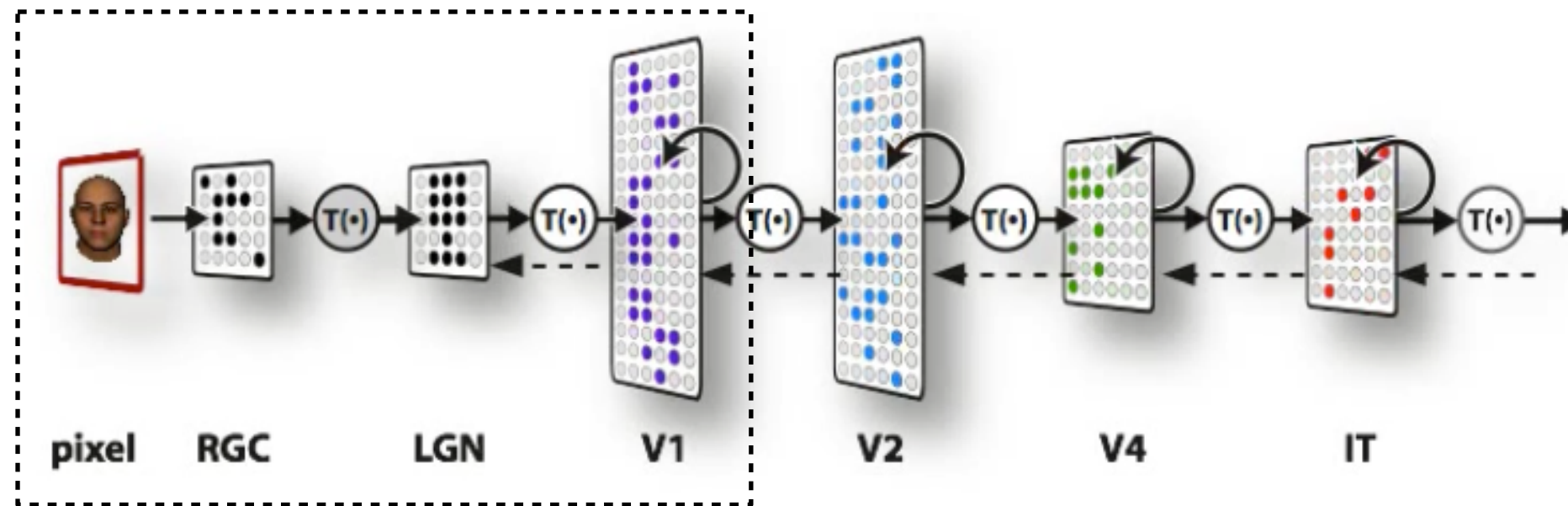Hubel and Wiesel (1965-1975) Lecun (2004), Carandini et. al (2005), Lennie & Movshon (2005) , DiCarlo (2012)

# Hierarchical Convolutional Neural Networks

▶ Individual layers of neurally-plausible **basic operations**

▶ Applied **convolutionally** — same at all locations:  approx. retinopy



*Operations in Linear-Nonlinear Layer*

$\otimes \Phi_1$
$\otimes \Phi_2$
...
$\otimes \Phi_k$
*Filter*

*Threshold*   *Pool*   *Normalize*

# Hierarchical Convolutional Neural Networks

▸ Individual layers of neurally-plausible **basic operations**

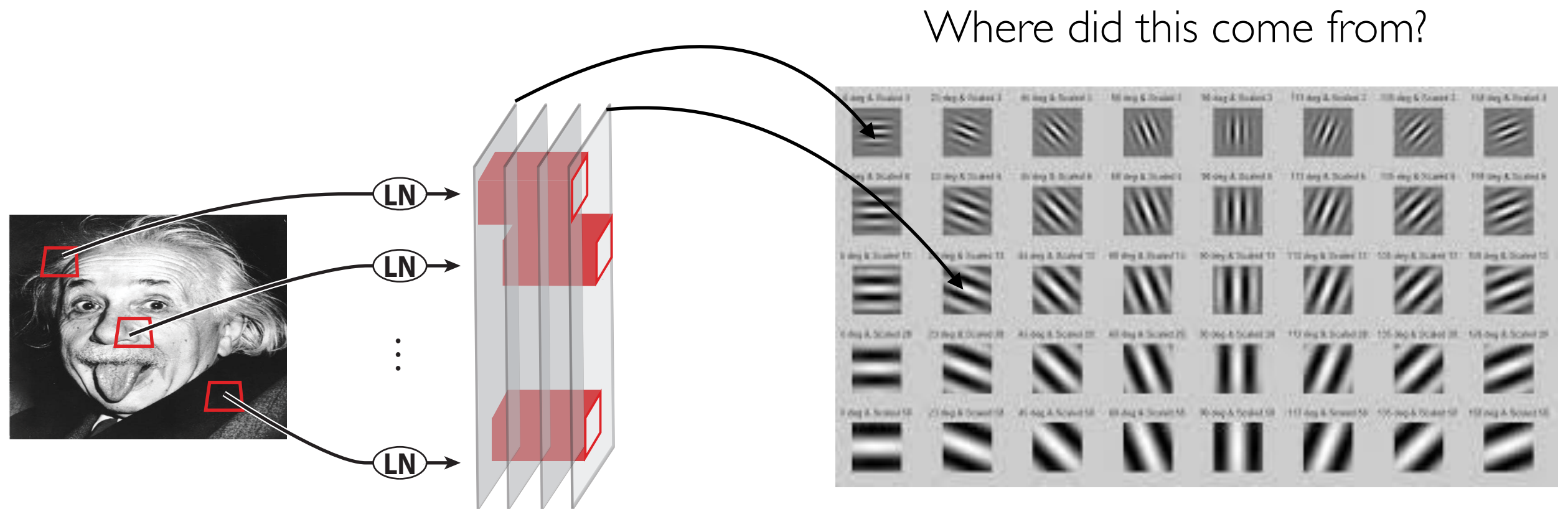▸ Applied **convolutionally** — same at all locations: approx. retinopy



image-input

image-like output

*Operations in Linear-Nonlinear Layer*

$\otimes \Phi_1$
$\otimes \Phi_2$
...
$\otimes \Phi_k$
*Filter*

*Threshold*    *Pool*    *Normalize*

▸ Individual layers of neurally-plausible **basic operations**

▸ Applied **convolutionally** — same at all locations: approx. retinopy



image-input

LN

LN

LN

one slice of output
for each filter pattern

*Operations in Linear-Nonlinear Layer*

$\otimes \Phi_1$
$\otimes \Phi_2$
...
$\otimes \Phi_k$
*Filter*

*Threshold*   *Pool*   *Normalize*

Lower areas, (RGC, LGN, V1) have been reasonably captured by single-layer convolutional model: ~50% of variance explained. Carandini et. al (2005), Lennie & Movshon (2005)
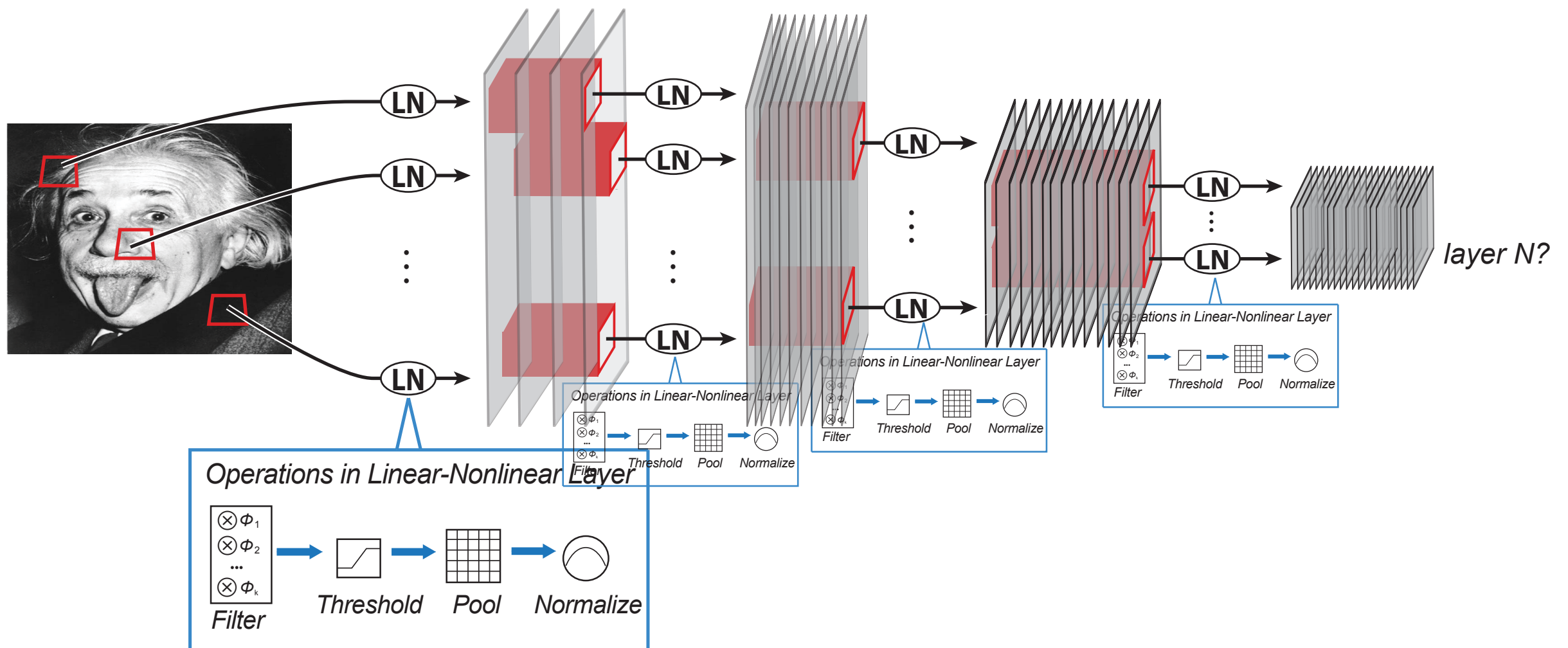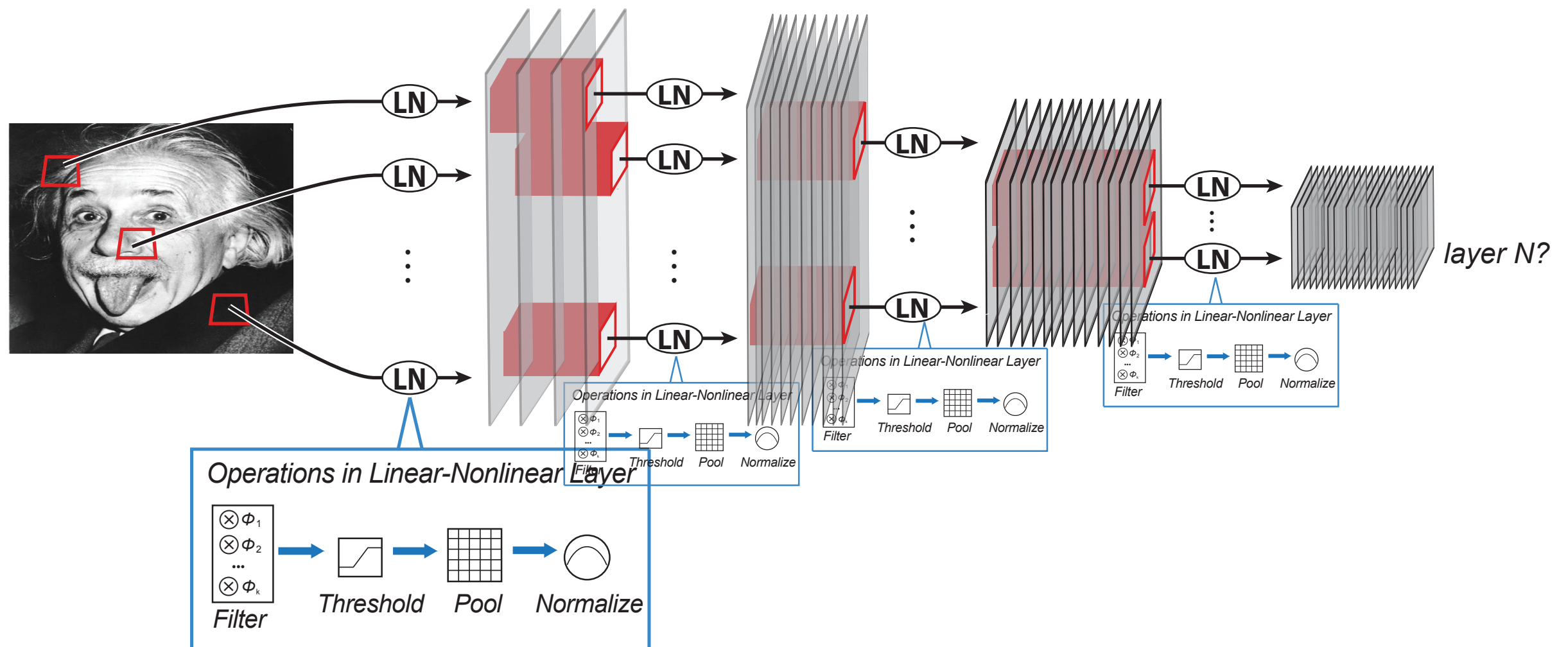
# Hierarchical Convolutional Neural Networks

Where did this come from?



(1) "Hubel and Wiesel's Intuition"
   ~1970s and formalized later

→ e.g. there is a "fixed basis set" that just "makes sense" if we're smart enough
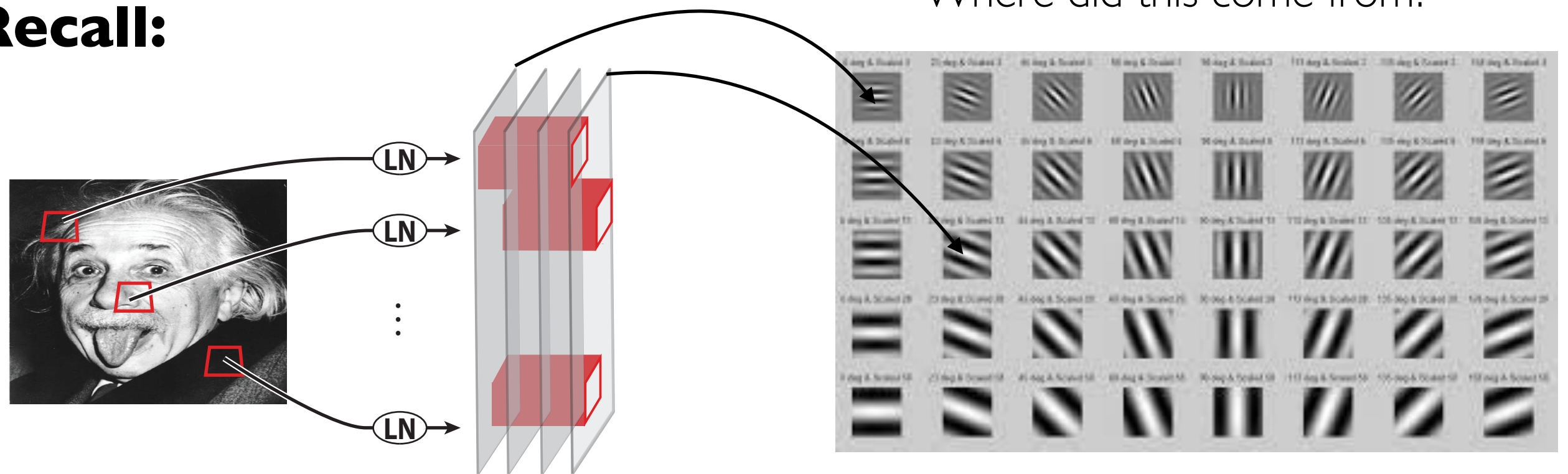
(2) Sparse Coding Foldiak, Olshausen, mid 1990s

→neurons have to represent their environment, as efficiently as possible

# Hierarchical Convolutional Neural Networks

Lower areas, (RGC, LGN, V1) have been reasonably captured by single-layer models: ~50% of variance explained. Carandini et. al (2005), Lennie & Movshon (2005)



Push up the ventral stream?

HCNNs

Huge number of parameters consistent with HCNN concept.



*layer N?*

Operations in Linear-Nonlinear Layer

$\otimes \Phi_1$
$\otimes \Phi_2$
...
$\otimes \Phi_k$

Filter    Threshold    Pool    Normalize

Huge number of parameters consistent with HCNN concept.



*Operations in Linear-Nonlinear Layer*

Filter — Threshold — Pool — Normalize

*layer N?*

i. **architectural** params:  (# layers, # filters, receptive field sizes, &c) — "network structure"

Huge number of parameters consistent with HCNN concept.



i. **architectural** params:  (# layers, # filters, receptive field sizes, &c) — "network structure"

ii. **filter** parameters:  continuous valued pattern templates — "network contents"

Huge number of parameters consistent with HCNN concept.



*Operations in Linear-Nonlinear Layer*

Filter    Threshold    Pool    Normalize

*layer N?*

i. **architectural** params: (# layers, # filters, receptive field sizes, &c) — "network structure"

ii. **filter** parameters: continuous valued pattern templates — "network contents"

*Q: How to discover the "right" parameters to understand real cortex?*

**Recall:**

Where did this come from?



(1) "Hubel and Wiesel's Intuition"
~1970s and formalized later

→ e.g. there is a "fixed basis set" that just "makes sense" if we're smart enough

(2) Sparse Coding Foldiak, Olshausen, mid 1990s

→ neurons have to represent their environment, as efficiently as possible

**Recall:**

Where did this come from?



(1) "Hubel and Wiesel's Intuition"
~1970s and formalized later

→ e.g. there is a "fixed basis set" that just "makes sense" if we're smart enough

## REALLY HARD TO GENERALIZE TO MULTI-LAYER NETWORKS

(2) Sparse Coding Foldiak, Olshausen, mid 1990s

→ neurons have to represent their environment, as efficiently as possible

Huge number of parameters consistent with HCNN concept.



*layer N?*

Operations in Linear-Nonlinear Layer

Filter   Threshold   Pool   Normalize

*Obvious alternative strategy: fit parameters to neural data.*

V4

10mm

IT

□ = Array

# Neural Fitting Strategy?

Huge number of parameters consistent with HCNN concept.

*layer N?*

*Operations in Linear-Nonlinear Layer*

Filter → Threshold → Pool → Normalize

…not enough neural data to constrain model class. Gallant (2007); Rust & Movshon (2006)

V4

IT

10mm

□ = Array

# Neural Fitting Strategy?

Huge number of parameters consistent with HCNN concept.



...not enough neural data to constrain model class. Gallant (2007); Rust & Movshon (2006)

Overfitting.

# Optimize for Performance, Test Against Neurons

Visual Recognition Task

Step 1: Optimize for Task

Spatial Convolution over Image Input

layer 1

layer 2

layer 3

layer 4

# Optimize for Performance, Test Against Neurons

Visual Recognition Task

Step 1: Optimize for Task

Spatial Convolution over Image Input

LN

layer 1

layer 2

layer 3

layer 4

Step 2: Compare to Neural Data

100ms Visual Presentation

V1

V2

V4

IT

1. **Performance**: accuracy on a <u>challenging, high-variation</u> visual object categorization task.

2. **Neural predictivity**: the ability of model to predict each individual neural site's activity.

1. **Performance**:  accuracy on a <u>challenging, high-variation</u>* visual object categorization task.

2. **Neural predictivity**: the ability of model to predict each individual neural site's activity.

**<u>*challenging for neural network engineers, not the animal</u>**

1. **Performance**:  accuracy on a <u>challenging, high-variation</u>* visual object categorization task.

2. **Neural predictivity**: the ability of model to predict each individual neural site's activity.

Our hypothesis:  Performance (1)  →  neural predictivity (2).

**<u>*challenging for neural network engineers, not the animal</u>**

High-throughput experiments to directly test the relationship between performance and IT neural predictivity.

▸ Random selection of model parameters; measure performance and neural predictivity   Pinto et. al (2008, 2009)

*Yamins\* and Hong\* et. al. **PNAS** (2014)*

High-throughput experiments to directly test the relationship between neural predictivity and performance.

▸ Random selection of model parameters; measure performance and neural predictivity    Pinto et. al (2008, 2009)

▸ Optimize parameters for performance; measure neural predictivity. optimization techniques: Bergstra Yamins & Cox (2013)

performance-
optimized
r = 0.79 ± .05
*(n=2000)*

Random selection
Performance Optimized

*Yamins* and Hong* et. al.* **PNAS** *(2014)*

High-throughput experiments to directly test the relationship between neural predictivity and performance.

▶ Random selection of model parameters; measure performance and neural predictivity   Pinto et. al (2008, 2009)

▶ Optimize parameters for performance; measure neural predictivity  optimization techniques: Bergstra Yamins & Cox (2013)

▶ Optimize parameters for neural predictivity; measure performance

predictivity-optimized
r = 0.80 ± .04
*(n=2000)*

Random selection
Performance Optimized
IT-Predictivity Optimized

IT Predictivity

Performance

*Yamins* and Hong* et. al.* **PNAS** *(2014)*

Performance is a potentially very good driver of neural prediction.



r = 0.55 ± .08
r = 0.79 ± .05
r = 0.80 ± .04

Random selection
Performance Optimized
IT-Predictivity Optimized

*Yamins\* and Hong\* et. al.* **PNAS** *(2014)*

But, not doing that well.   Really want to be here:



IT Predictivity

V1-like

V2-like

HMAX

SIFT

PLOS09

Pixels

Random selection
Performance Optimized
IT-predictivity Optimized

Performance

i. **architectural** params: (# layers, # filters, receptive field sizes, &c) — "network structure"

→ Automated meta-parameter optimization in high-dimensional discrete parameter spaces
Bergstra Yamins & Cox (2013)

→ Ensembles of models chosen through modified boosting Yamins et. al (2013, 2014)

# Optimization Strategy

i. **architectural** params: (# layers, # filters, receptive field sizes, &c) — "network structure"

→ Automated meta-parameter optimization in high-dimensional discrete parameter spaces
Bergstra Yamins & Cox (2013)

→ Ensembles of models chosen through modified boosting  Yamins et. al (2013, 2014)

ii. **filter** parameters: continuous valued pattern templates — "network contents"

→ GPU-accelerated stochastic gradient descent  Pinto et. al., (2009), Krizhevsky et. al. (2012)

Gradient descent eq:
$$\frac{dp_a}{dt} = -\lambda(t) \cdot \langle \nabla_{p_a} L(x) \rangle_{x \in \mathcal{D}}$$

L = loss function
$\lambda$ = learning rate
D = dataset

In current practice:

L = loss computed from **large numbers of externally-provided object category labels.**

# Model Training Regimen

ImageNet (2012). Thousands of images in thousands of categories.

**train:** real photos

**train:** real photos

**test:** neural stimuli



generalize?

Basic
categorization

# Model Training Regimen

**train:** real photos



**test:** neural stimuli

generalize?

Basic categorization

removed categories of photos that
appeared in the test stimuli
(animals, boats, cars, chairs, faces, fruits, planes, tables)

# Model Training Regimen

**train:** real photos

**test:** neural stimuli



generalize?

Basic categorization

removed categories of photos that appeared in the test stimuli
(animals, boats, cars, chairs, faces, fruits, planes, tables)

→ Specific 4-layer model that achieved high recognition performance.

# Performance Generalization



Basic categorization

test performance

0.75

0.55

0.35

training time →

# IT Neurons Track Human Performance

V4 loses out at higher variation:

… but humans are much less affected.

… as is the IT neural population.

Basic categorization

Performance *(% correct)*

Pixels · SIFT · V1-like · V2-like · HMAX · PLOS09 · V4 NEURONS · IT NEURONS · HUMAN PERFORMANCE

Low Variation — Medium Variation — High Variation

V4-to-IT Gap

*Yamins\* and Hong\* et. al.* **PNAS** *(2014)*

At <u>high variation levels</u>, IT much better than V4 and existing models.

# Performance Comparison

At <u>high variation levels</u>, IT much better than V4 and existing models

Performance
*(% correct)*

100
80
60
40
20

Pixels
SIFT
V1-like
V2-like
HMAX
PLOS09
V4 NEURONS
IT NEURONS
HMO MODEL
HUMAN PERFORMANCE

Low Variation

Medium Variation

V4-to-IT

Gap

High Variation

New model comparable to IT / human performance levels.

**GOAL:  Predictive model of single-neuron responses throughout the ventral stream to arbitrary image stimuli.**



pixel   RGC         LGN         V1          V2          V4          IT

Get quantitative hypothesis
for the network
that generated
this data  —>

Img 1        Img 2              Img **5760**

Neuron 1
Neuron 2
Neuron 3

Neuron **296**

# Predicting IT Neural Responses



IT Site 150          IT Site 56          IT Site 42

Response Magnitude

Images sorted first by **category**, then **variation level**.

——— Neural data

——— Model prediction

# Key Underlying Principle

Yamins* and Hong* et. al. **PNAS** (2014)

$r = 0.87 \pm 0.15$

HMO

V2-like

SIFT

HMAX

PLOS09

V1-like

Pixels

○ = distinct model

IT Explained Variance (%)

Categorization Performance (balanced accuracy)

Captures low variation image response patterns ...

Layer I

Animals   Boats   Cars   Chairs   Faces   Fruits   Planes   Tables

Neural data

Model prediction

Layer
I

Animals    Boats    Cars    Chairs    Faces    Fruits    Planes    Tables

Neural data

Model prediction

… but fails to capture higher variation response patterns.

Layer **2**

Layer **1**

Animals  Boats  Cars  Chairs  Faces  Fruits  Planes  Tables

—— Neural data      —— Model prediction

Layer **3**

Layer **2**

Layer **1**

Animals    Boats    Cars    Chairs    Faces    Fruits    Planes    Tables

Neural data      Model prediction

Building tolerance while maintaining selectivity

Top Layer

Layer 3

Layer 2

Layer 1

Animals   Boats   Cars   Chairs   Faces   Fruits   Planes   Tables

# Predicting IT Neural Responses

# Predicting IT Neural Responses

# Predicting IT Neural Responses

Performance constraints

*Yamins\* and Hong\* et. al.* **PNAS** *(2014)*

# Predicting IT Neural Responses



Yamins* and Hong* et. al. **PNAS** (2014)

# Predicting IT Neural Responses

Performance constraints + architectural constraints → better neural prediction

IT Explained Variance (%)

**Ideal Observers:** Category, All Variables

**Control Models:** Pixels, V1-Like, SIFT, PLOS09, HMAX, V2-Like

**HMO Layers:** HMO L1, HMO L2, HMO L3, HMO Top

What about intermediate layers?

i. compare intermediate model layers to IT neural data

ii. compare all model layers to intermediate visual areas (V4)

V4 unit 60

# Predicting V4 Neural Responses



Top Layer

Layer **3**

Layer **2**

Layer **1**

Animals  Boats  Cars  Chairs  Faces  Fruits  Planes  Tables

Neural data

Model prediction

# Predicting V4 Neural Responses

*Yamins\* and Hong\* et. al.* **PNAS** *(2014)*

Investigating fits as a function of model layer:



Yamins* and Hong* et. al. **PNAS** (2014)

IT fit increases at each layer.  In contrast, V4 fit peaks and then goes down.

Model output at lowest layer resembles Gabor wavelets:



Layer 1 Filters

In submission: model lowest layer is best explanation of imaging data in V1. (*with Darren Seibert and Justin Gardner*)

Complement standard "from below" approach …



pixel   RGC        LGN        V1        V2        V4        IT

Complement standard "from below" approach …



pixel    RGC         LGN         V1          V2          V4          IT

Complement standard "from below" approach …



pixel   RGC      LGN      V1      V2      V4      IT

Complement standard "from below" approach …



pixel   RGC   LGN   V1   V2   V4   IT

Complement standard "from below" approach ... with behavioral constraints

Complement standard "from below" approach … with behavioral constraints



pixel    RGC         LGN         V1          V2          V4          IT

Complement standard "from below" approach … with behavioral constraints



pixel    RGC    LGN    V1    V2    V4    IT

Complement standard "from below" approach … with behavioral constraints



pixel   RGC   LGN   V1   V2   V4   IT

Similar ideas and results:
Khaligh-Razavi & Kriegeskorte (2014),
Guclu & Van Gerven (2015), Cichy & Oliva (2015)

*plane*  *Category*

*f16*  *Identity*

*Position*

*Size*

*Aspect Ratio and Angle*

# Beyond categorization

We can quickly assess the scene as a whole.



*plane* — Category

*f16* — Identity

Bounding Box

X and Y Axis Position

Aspect Ratio

Major Axis Length

Major Axis Angle

Perimeter

2-D Retinal Area

3-D Object Scale

*rz* *rx* *ry* — Pose in each axis

# Where and how are all these properties coded neurally?

"Standard word model" predicts: **not at the top of the ventral stream**.

Aggregation over identity-preserving transformations, e.g. translation.

# Beyond categorization

"Standard word model" predicts: **not at the top of the ventral stream**.

Aggregation over identity-preserving transformations, e.g. translation.



Receptive Field Size ↑

Category Invariance ↑

# Beyond categorization

"Standard word model" predicts: **not at the top of the ventral stream**.

Aggregation over identity-preserving transformations, e.g. translation.



Receptive Field Size ↑

Category Invariance ↑

(e.g.) Position Sensitivity ↓

categorization

**IT**

pose?
**V4**

**V2**

**V1**

position / size estimation

# Where and how are all these properties coded neurally?



dorsal stream?

earlier visual areas?

V1

V4

V2

IT

Category
plane

Identity
f16

Position

Size

Bounding Box

Aspect and Angle

Pose

# Beyond categorization

Unexpected observation:

Training on categorization task $\longrightarrow$ *Increased* performance on position estimation task.

*even though the goal was to become INVARIANT to position*

# Beyond categorization

Category optimization → improved performance on non-categorical tasks.



*Hong\*, Yamins\*, Majaj & DiCarlo.* **Nat. Neuro.** *(2016)*

# Beyond categorization

Unexpected observation #2:



*Hong\*, Yamins\*, Majaj & DiCarlo.* **Nat. Neuro.** *(2016)*

# Beyond categorization

Unexpected observation #2:



*Increased* performance on
position estimation task
at each model layer.



Categorization

Test Performance

Layer 1  Layer 2  Layer 3  Layer 4  Layer 5  Layer 6

Model Layers

X-Axis Position

*Hong\*, Yamins\*, Majaj & DiCarlo.* **Nat. Neuro.** *(2016)*

For all tasks of visual interest we could measure in our test dataset:



Performance on non-categorical tasks increases at each layer.

What do the data say?

category: plane
identity: f16

Categorization

Identification

Pix V1 V4 IT

Hong*, Yamins*, Majaj & DiCarlo. **Nat. Neuro.** (2016)

IT cortex

V4 cortex

V1-like model

pixel control

# Population Decoding



category: plane
identity: f16

Categorization

Identification

X-axis Position

Y-axis Position

Hong*, Yamins*, Majaj & DiCarlo. **Nat. Neuro.** (2016)

IT cortex

V4 cortex

V1-like model

pixel control

Population Decoding

IT > V4, V1 for all tasks    V4 > V1 for most tasks

*Hong\*, Yamins\*, Majaj & DiCarlo. Nat. Neuro. (2016)*

IT cortex    V4 cortex
V1-like model    pixel control

"Standard" receptive field-mapping stimuli w/ position and orientation variation:

X-position

Y-position

Orientation

# Population Decoding

## V1 > V4, IT   for "standard" tasks



Hong*, Yamins*, Majaj & DiCarlo. **Nat. Neuro.** (2016)

IT cortex — V4 cortex — V1-like model — pixel control

Click where the **boat** was!

**6 learning trial(s) left** after this.

# Monkey Neurons vs Humans

$$\text{performance} \sim k * \log(N)$$



Basic Categorization — Subordinate Identification

*Hong\*, Yamins\*, Majaj & DiCarlo.* **Nat. Neuro.** *(2016)*

# Monkey Neurons vs Humans



Hong*, Yamins*, Majaj & DiCarlo. **Nat. Neuro.** (2016)

Somewhat newish ideas about IT?

State of knowledge from previous studies . . .

Multiple hypotheses consistent with the existing data . . .

Population Decode Performance *(relative to human performance)*

Depth Along Ventral Stream *(increasing receptive field size →)*

Categorization

Orthogonal Properties

H1

H2

H3

H3: Information preservation?

Somewhat newish ideas about IT?

Somewhat newish ideas about IT?

H4: Simultaneous build-up of encoding

1. IT is *NOT* invariant. Strict generalization of simple-to-complex cells: **no**.

2. "Lower-level" properties are not that low-level — at least, with complex objects and backgrounds.

**3**. Categorization and non-categorical properties "go together" — *not* just that "not all (e.g.) position information is lost" (MacEvoy 2013, DiCarlo 2003)

Provides support to a hypothesis for what IT does:

"Inverting the generative model of the scene"

Alex Kell

Sam Norman-Haignere

Josh McDermott

How are circuits making sense of complex sound patterns?

# Core / Belt / Parabelt Structure



Core area

Belt area

A1

P

A

Parabelt area

*monkey
*

*Tramo et. al, Curr. Opin. Neuro. (1999)*

Core area

Belt area

A1

P

A

*Tramo et. al, Curr. Opin. Neuro. (1999)*

Parabelt area

\*monkey
\*

# Core / Belt / Parabelt Structure



Core area

Belt area

A1

P

A

Tramo et. al, Curr. Opin. Neuro. (1999)

Parabelt area

*monkey
*

Spatiotemporal filtering? *Shamma, 2005*



Core area

Belt area

A1

P

A

Parabelt area

*Tramo et. al, Curr. Opin. Neuro. (1999)*

*__monkey__

*

Spatiotemporal filtering? *Shamma, 2005*

Core area

??? Belt area
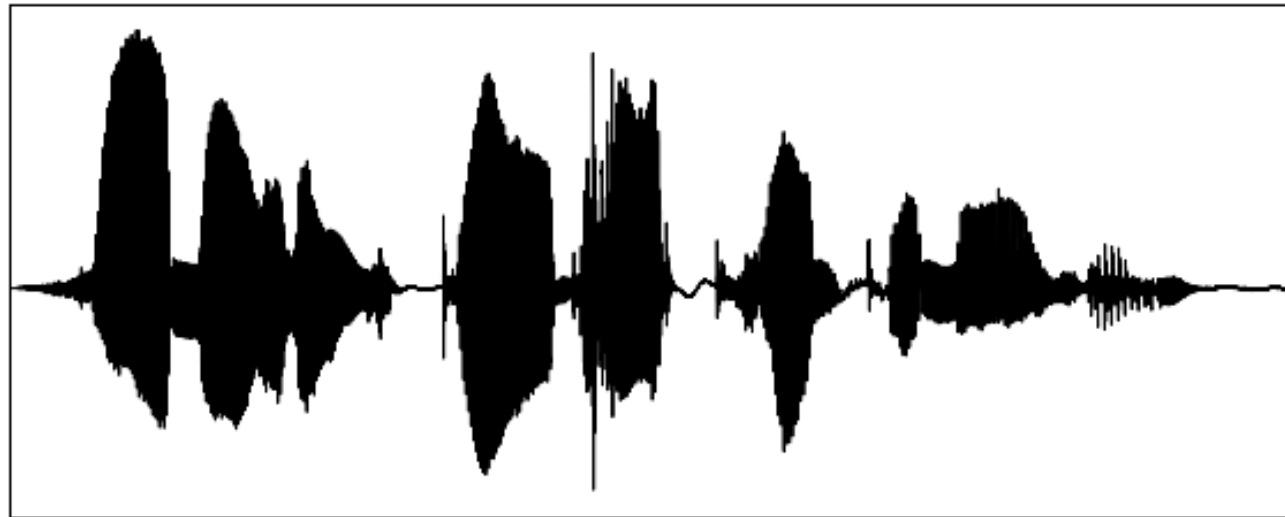


*Tramo et. al, Curr. Opin. Neuro. (1999)*

*monkey

*

Parabelt area

???

Spatiotemporal filtering? *Shamma, 2005*

**???**

Core area

Belt area

*Tramo et. al, Curr. Opin. Neuro. (1999)*

A1

P

A

*monkey

*

Parabelt area

**???**

Our goal: use computational models to help deepen understanding of non-primary areas.
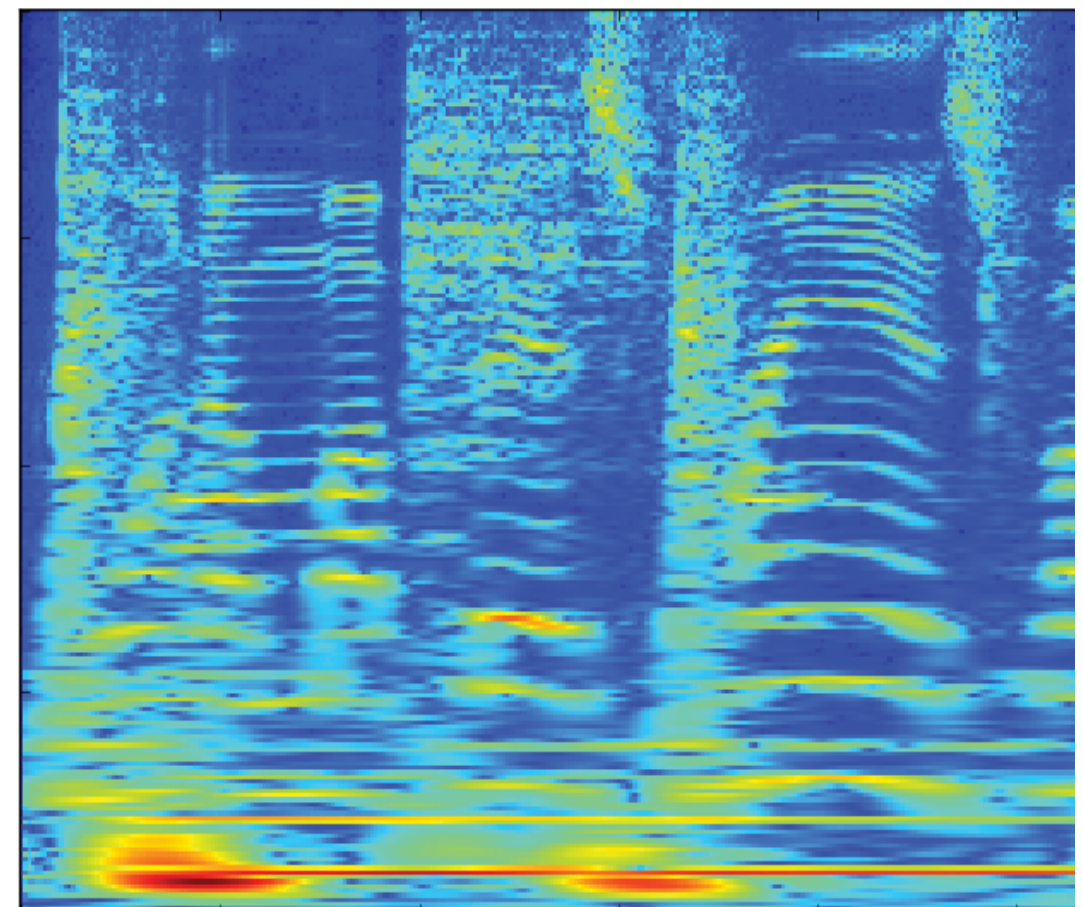
# Convolutional Neural Networks
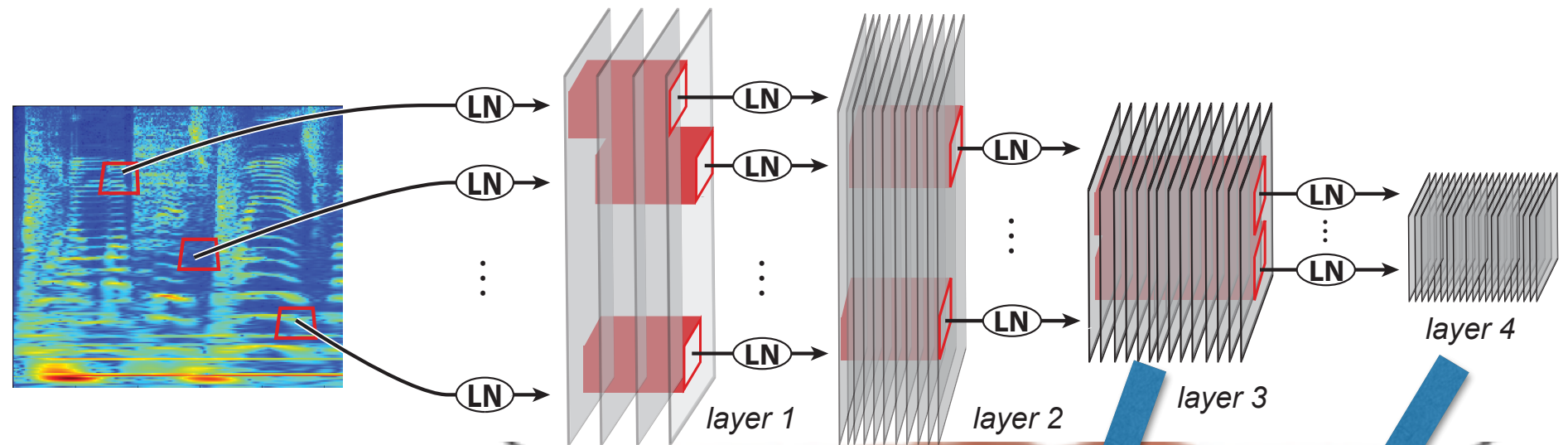
Waveform representation



Time →

Cochleagram representation



Frequency →
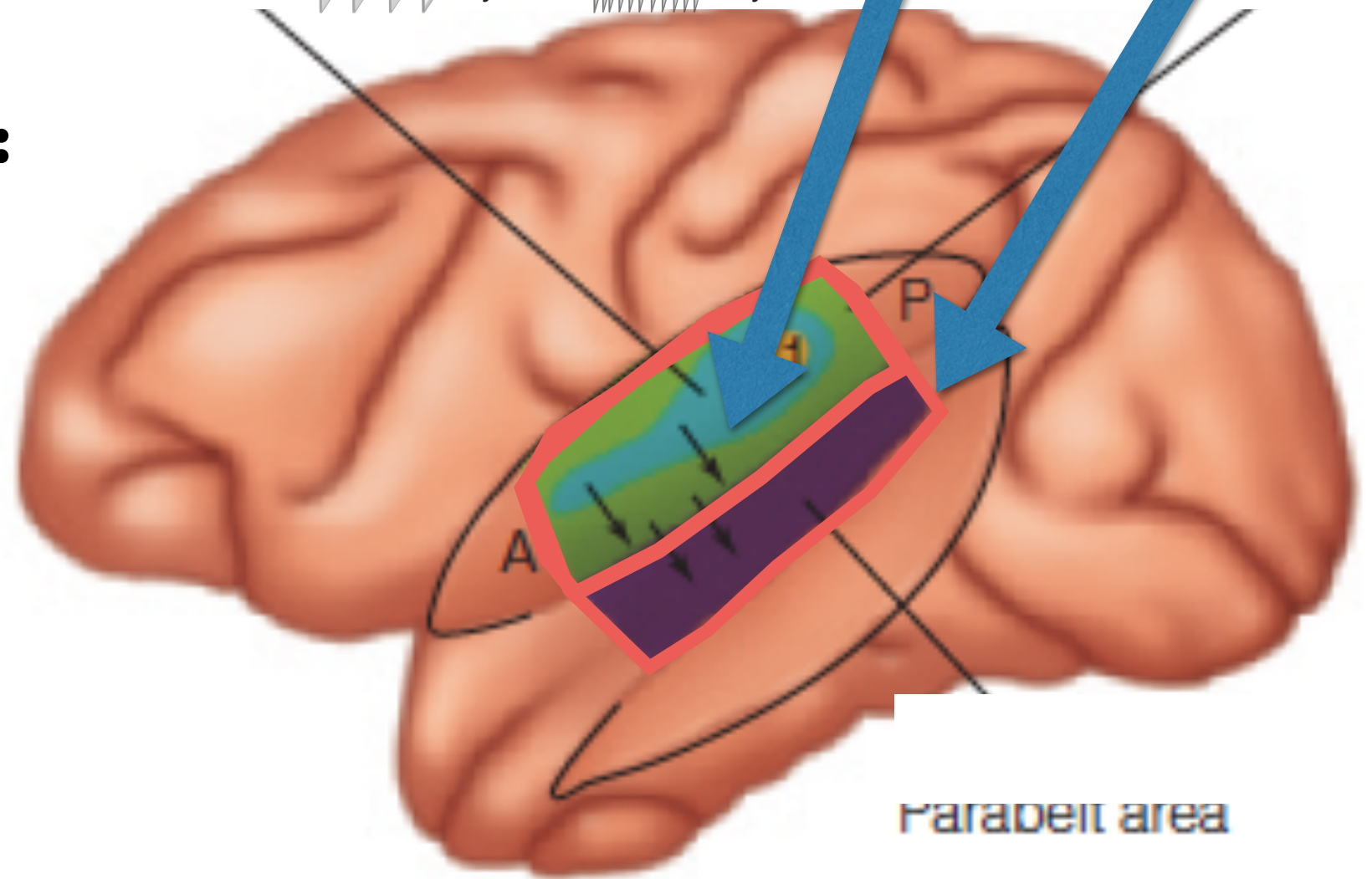
Time →

Coarse model of the cochlea

**Task-Driven Modeling:**

1. Optimize for performance on a challenging auditory task, fixing parameters

2. Compare to neural data.

Apply to auditory tasks, where the regions themselves are less well known.

**600-way** word-recognition task assembled by:

- Recordings from standard speech recognition databases (TIMIT, WSJ) with words spoken at least 20 times

- Combined with significant background noise

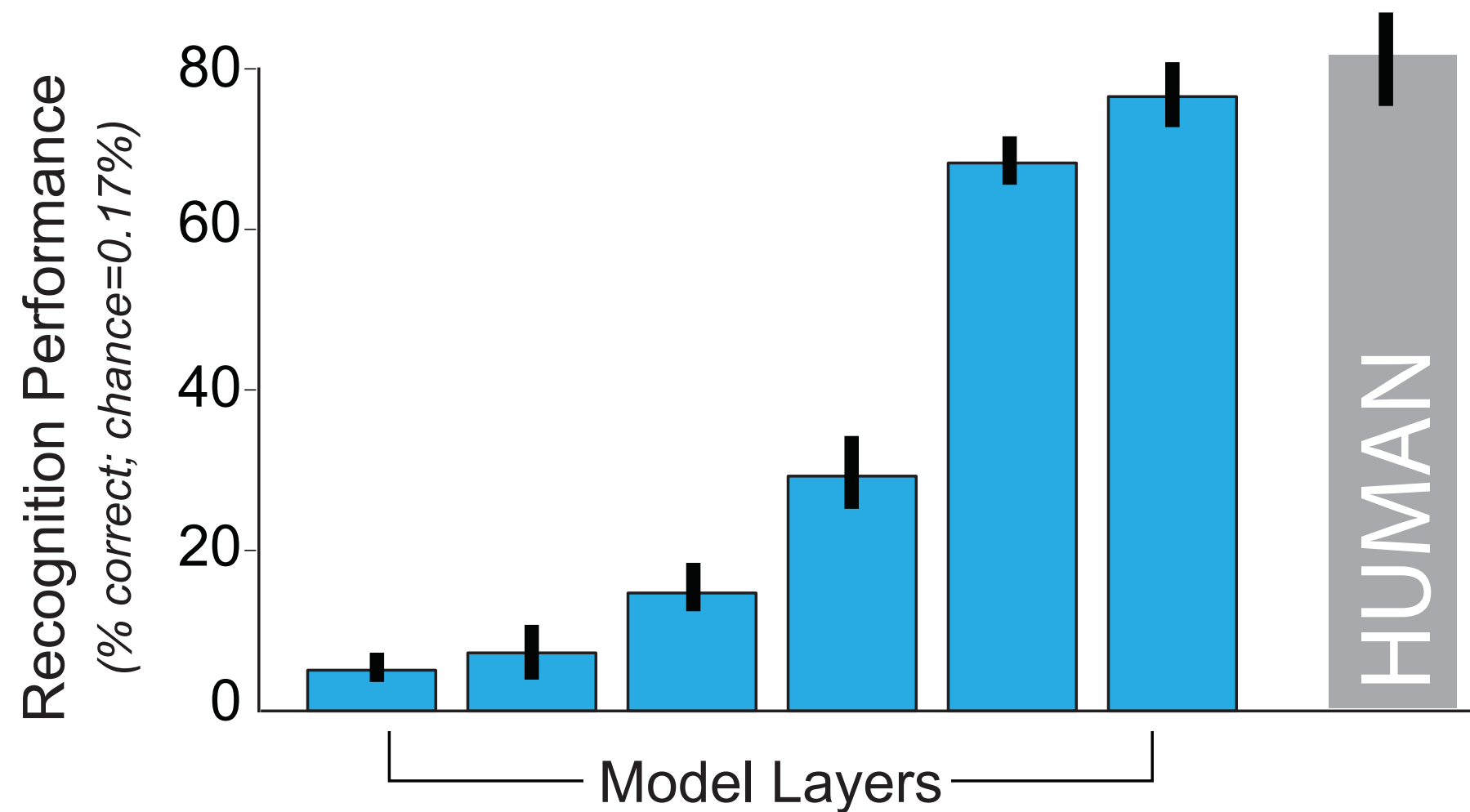▸ auditory scenes

▸ speech babble

▸ music clips

*"She **had** your*      *'had'*

*dark **suit** in*      *'suit'*

*greasy **wash** water*      *'wash'*

*all **year** … "*      *'year'*

**600-way** word-recognition task assembled by:

- Recordings from standard speech recognition databases (TIMIT, WSJ) with words spoken at least 20 times

- Combined with significant background noise

▶ auditory scenes

▶ speech babble

▶ music clips



Backgrounds → humans not close to ceiling.

# Performance Results

Performance on 600-way word-recognition task



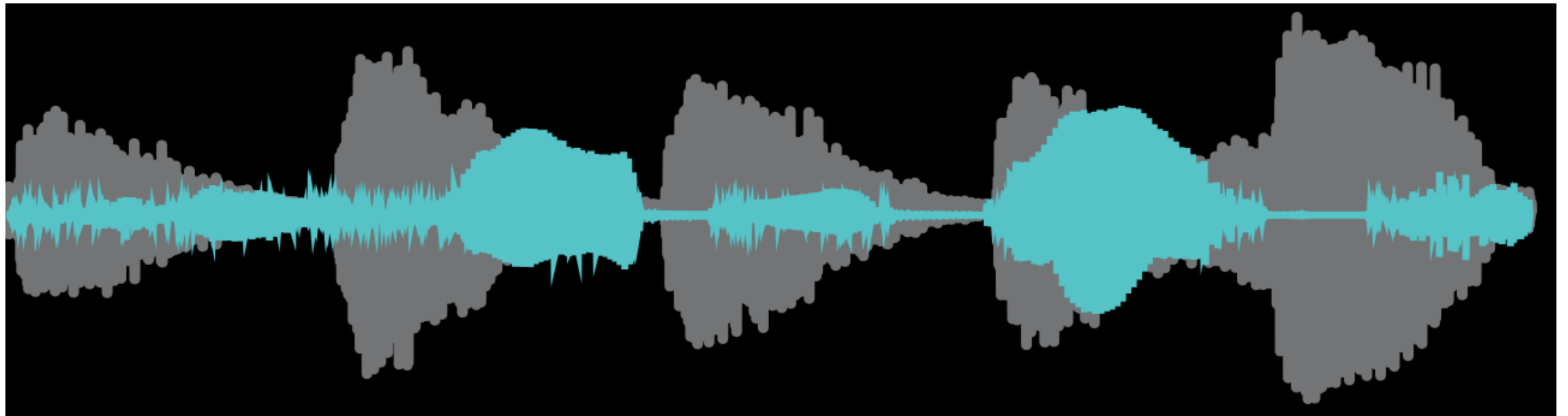… for model, measured on held-out data with novel speakers and auditory background noise.

# Behavioral comparison:
# CNN & humans on same task



Word recognition in complex backgrounds

# Behavioral comparison:
# CNN & humans on same task



## Word recognition in complex backgrounds

21 conditions:
dry
+
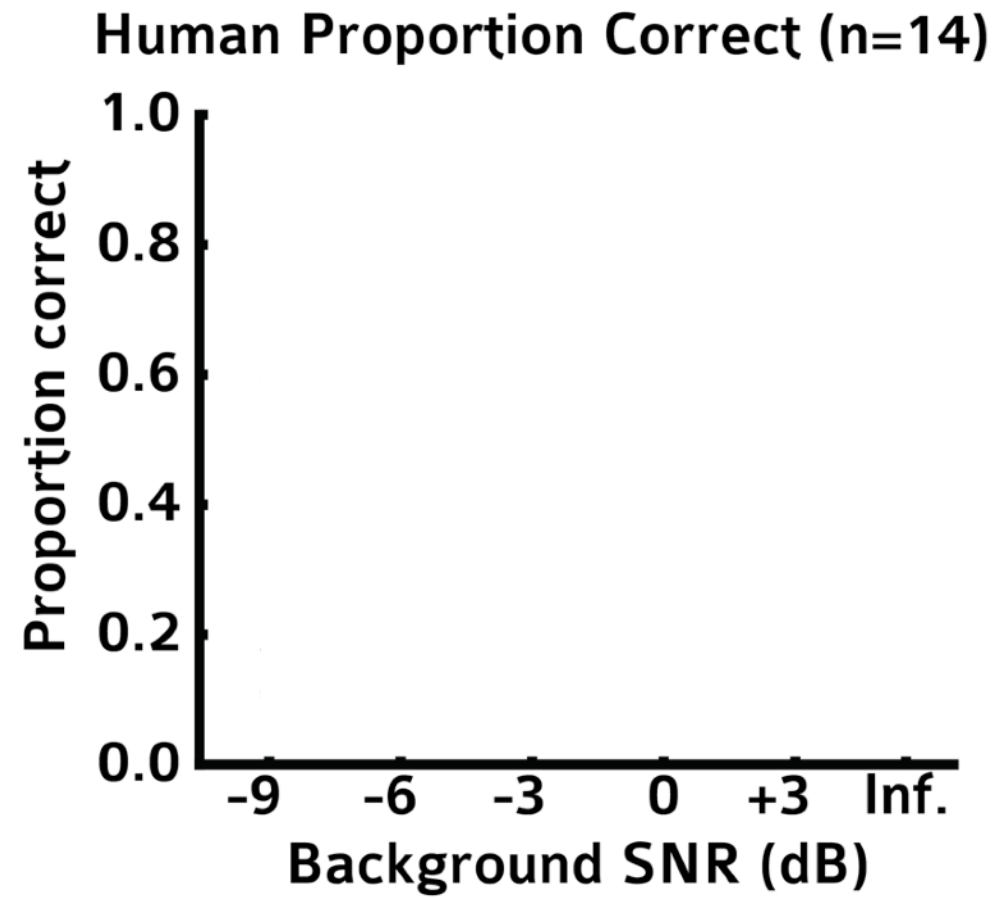4 different background types at 5 SNR levels:
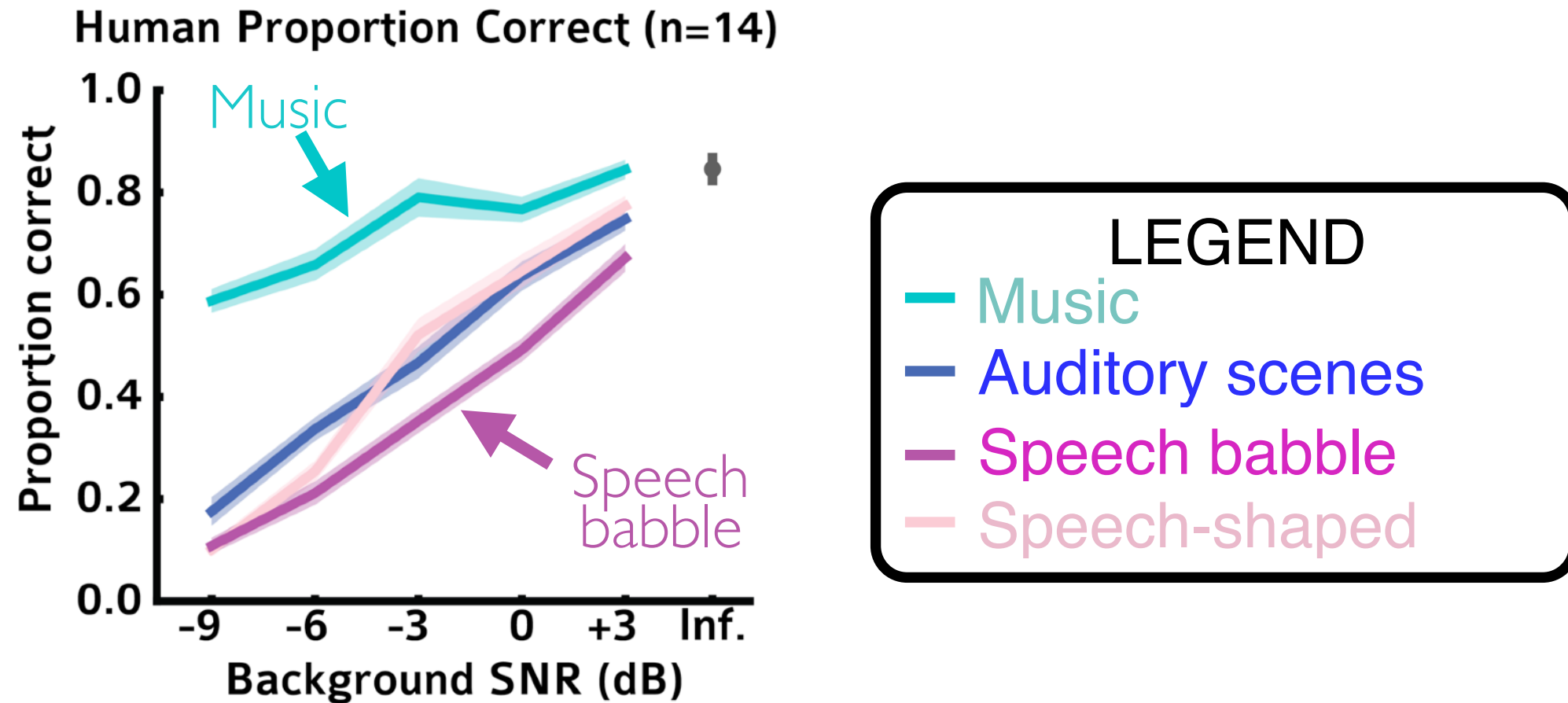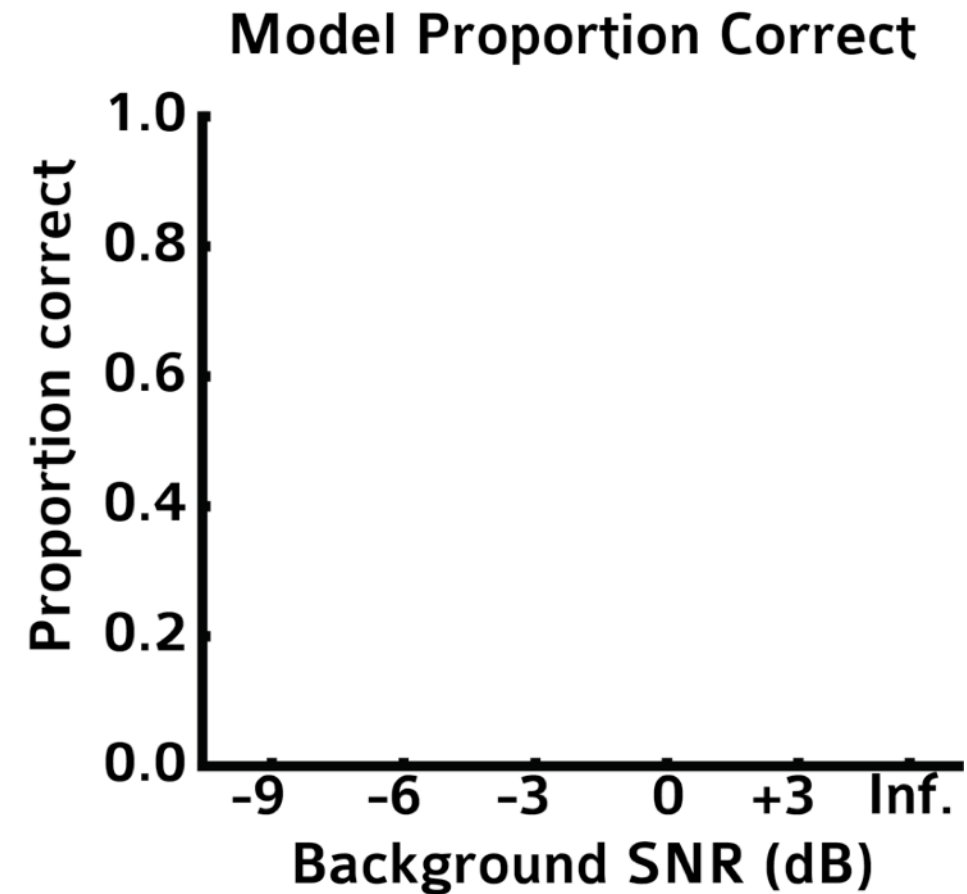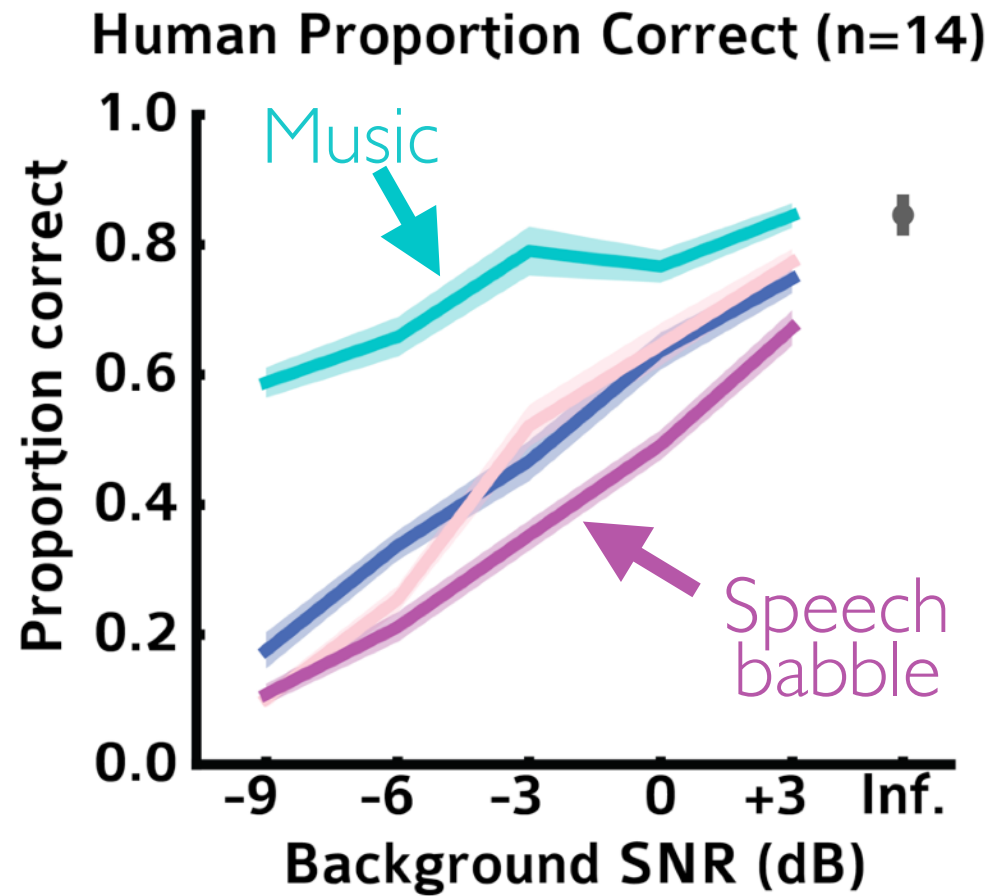
Auditory scenes          Speech babble
Music          Speech-shaped noise

# Behavioral comparison:
# CNN & humans on same task



Word recognition in complex backgrounds

21 conditions:
dry
+
4 different background types at 5 SNR levels:

600
AFC

Auditory scenes    Speech babble

Music    Speech-shaped noise

# Behavioral comparison: CNN & humans on same task



**Human Proportion Correct (n=14)**

Y-axis: Proportion correct (0.0, 0.2, 0.4, 0.6, 0.8, 1.0)

X-axis: Background SNR (dB) (-9, -6, -3, 0, +3, Inf.)

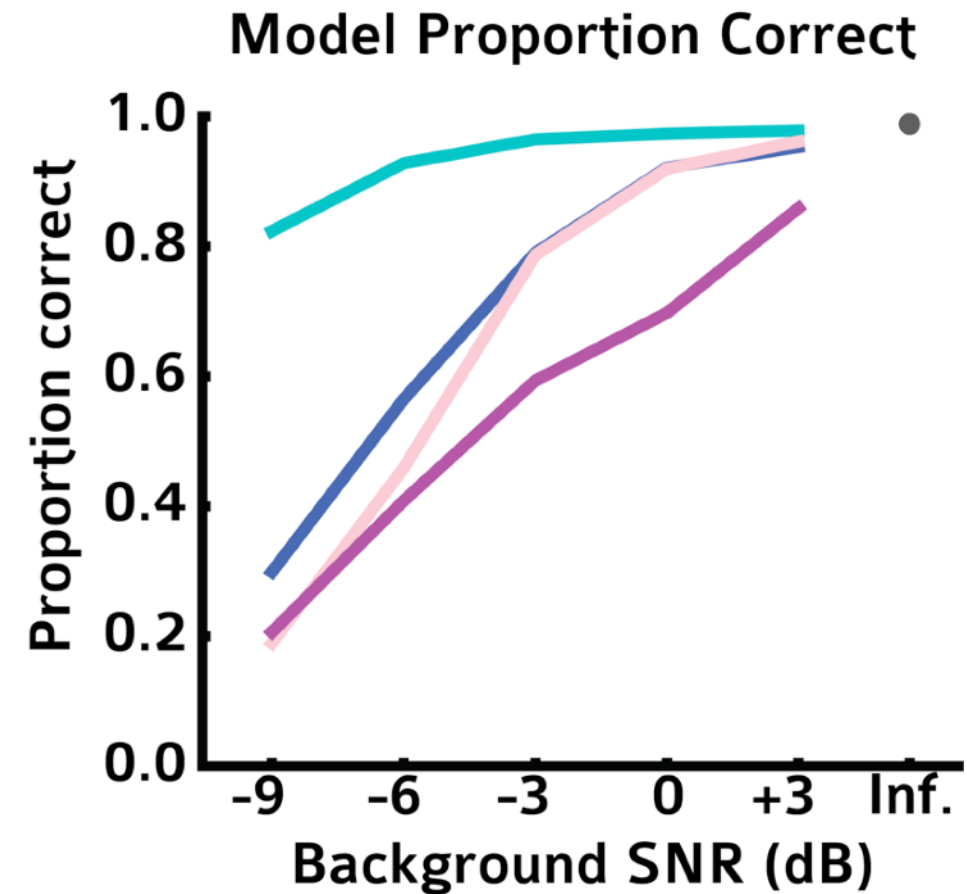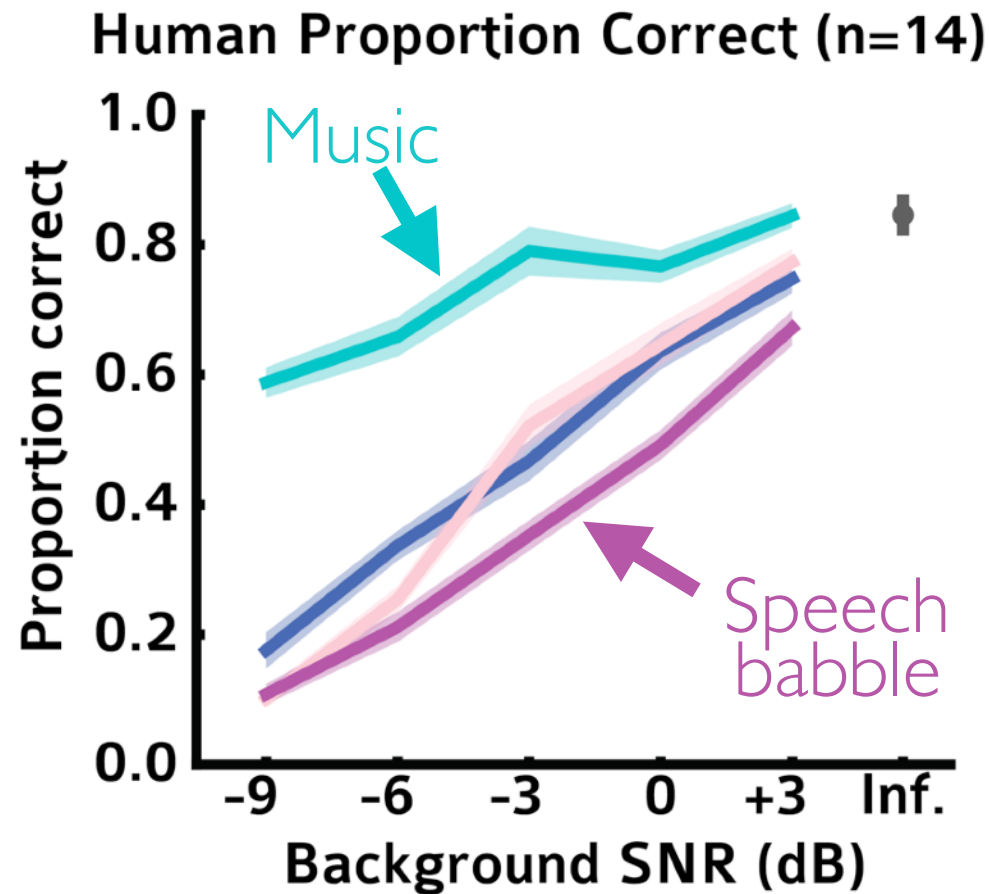# Behavioral comparison: CNN & humans on same task



Human Proportion Correct (n=14)

LEGEND
— Music
— Auditory scenes
— Speech babble
— Speech-shaped
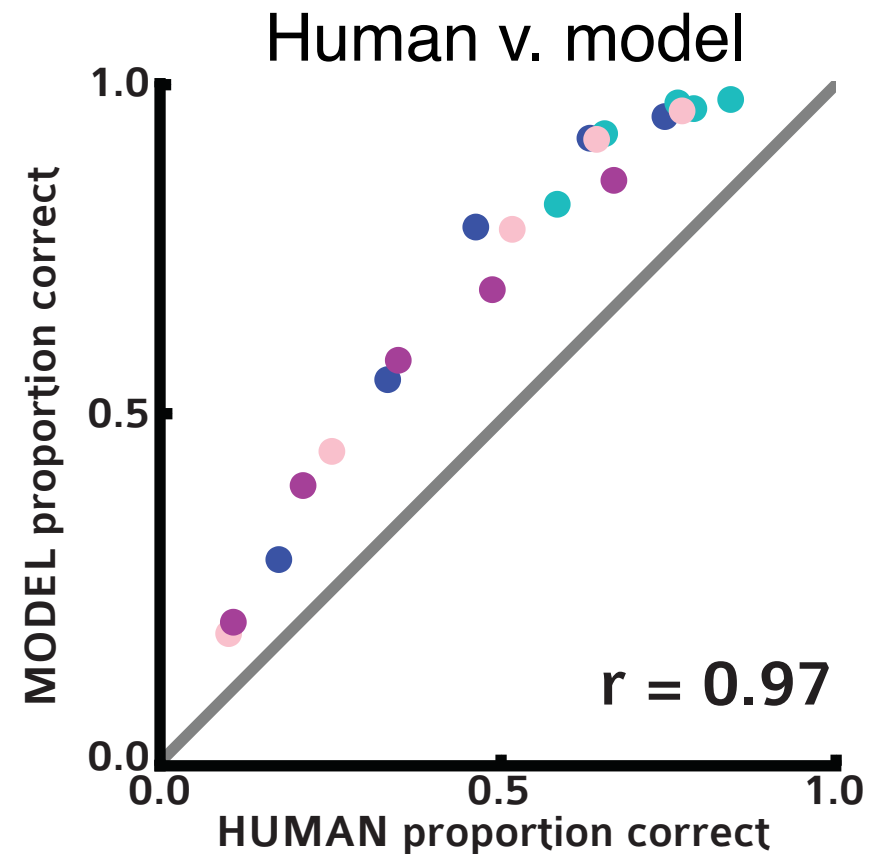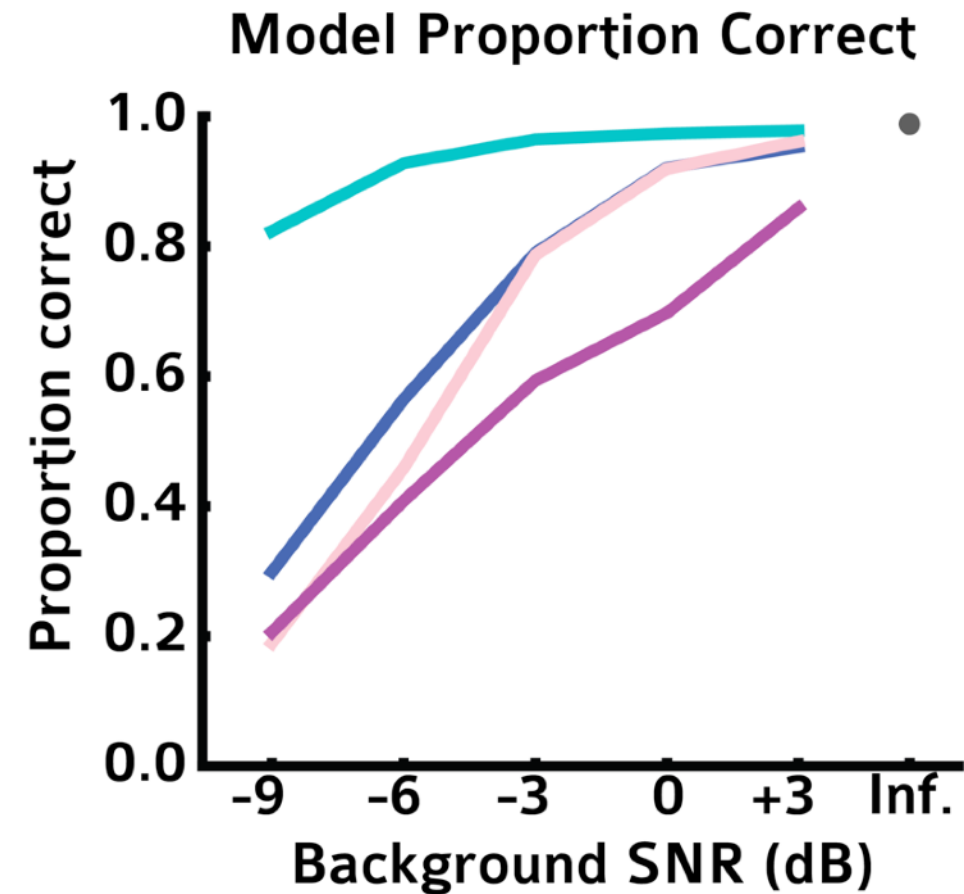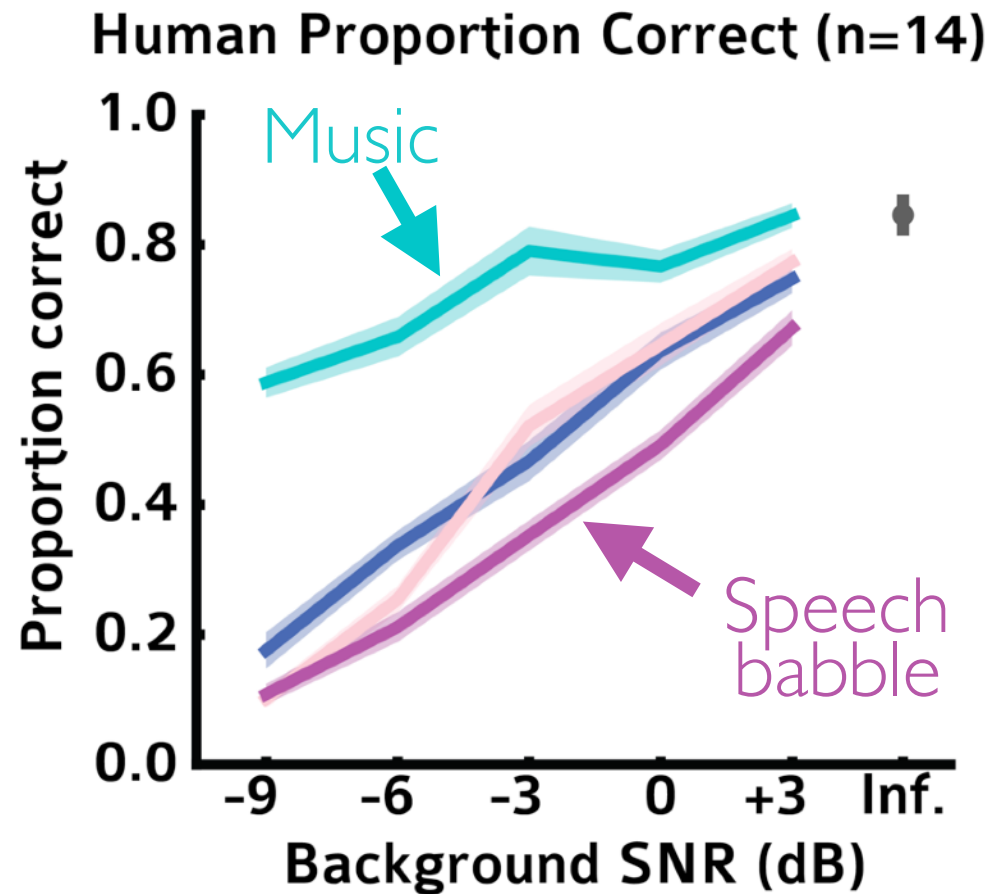
# Behavioral comparison: CNN & humans on same task

# Behavioral comparison: CNN & humans on same task



**Human Proportion Correct (n=14)**

Music

Speech babble

Proportion correct

Background SNR (dB)

-9 -6 -3 0 +3 Inf.

**Model Proportion Correct**

Proportion correct

Background SNR (dB)

-9 -6 -3 0 +3 Inf.

LEGEND
— Music
— Auditory scenes
— Speech babble
— Speech-shaped

# Behavioral comparison: CNN & humans on same task

# Behavioral comparison: CNN & humans on same task



**Human Proportion Correct (n=14)**

Music

Speech babble

Background SNR (dB)

**Model Proportion Correct**

Background SNR (dB)

**Human v. model**

r = 0.97

MODEL proportion correct

HUMAN proportion correct

LEGEND
- Music
- Auditory scenes
- Speech babble
- Speech-shaped

# Behavioral comparison: CNN & humans on same task

# Does distortion in a periphery-like representation explain pattern of performance?

## Measure physical distortion of background noise



Dry        Wet        |Dry - Wet|

Freq.

Time

# Distortion and pattern of performance.



**Human Proportion Correct (n=14)**

Music

Speech babble

Proportion correct

Background SNR (dB): -9, -6, -3, 0, +3, Inf.

Cochleagram distortion by condition
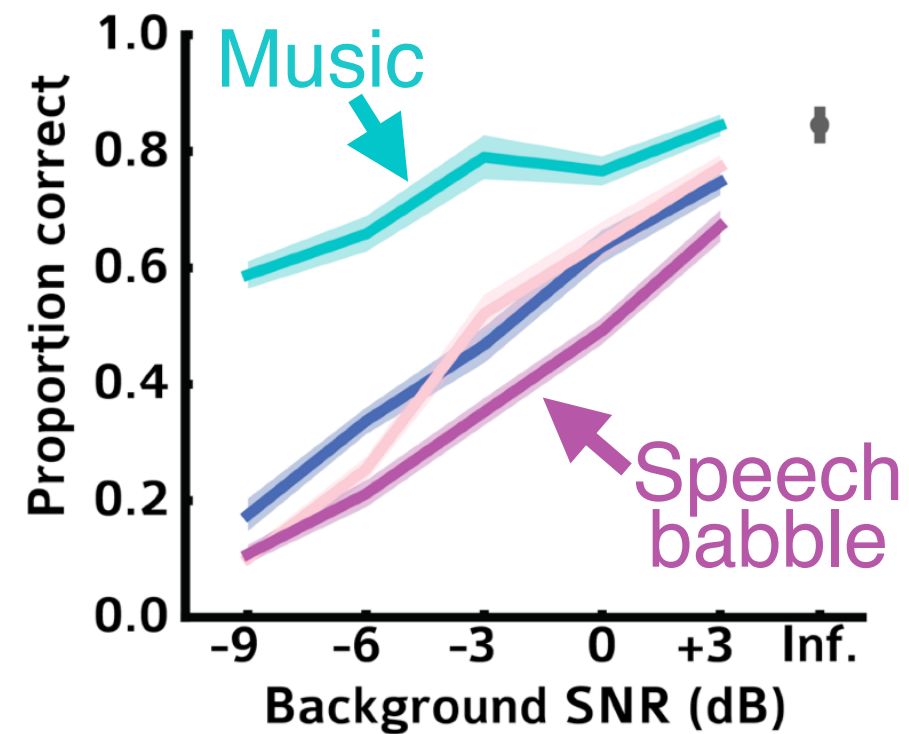
Less distortion

More distortion

Background SNR (dB): -9, -6, -3, 0, +3

LEGEND

● Music  ● Auditory scenes  ● Speech babble  ● Speech-shaped noise

# Distortion and pattern of performance.

# Distortion and pattern of performance.
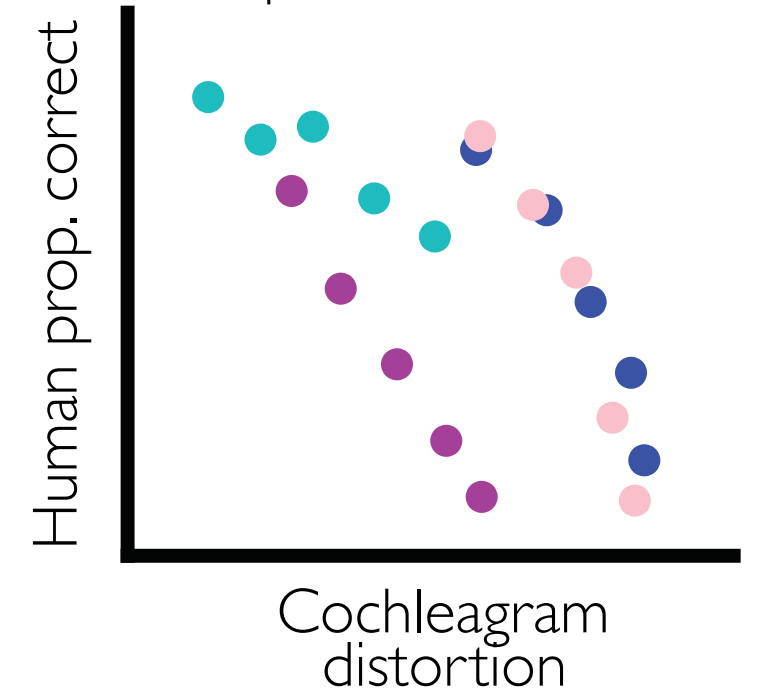
# Distortion and pattern of performance.



**Human Proportion Correct (n=14)**

Music

Speech babble

Proportion correct

Background SNR (dB)

-9  -6  -3  0  +3  Inf.

Cochleagram distortion by condition

Less distortion

More distortion

Background SNR (dB)

-9  -6  -3  0  +3

Cochleagram distortion v. human performance
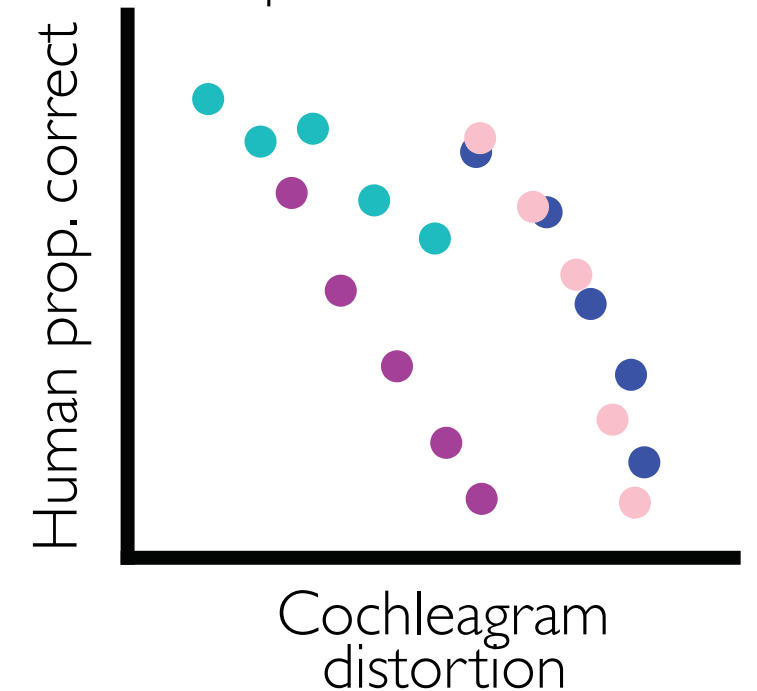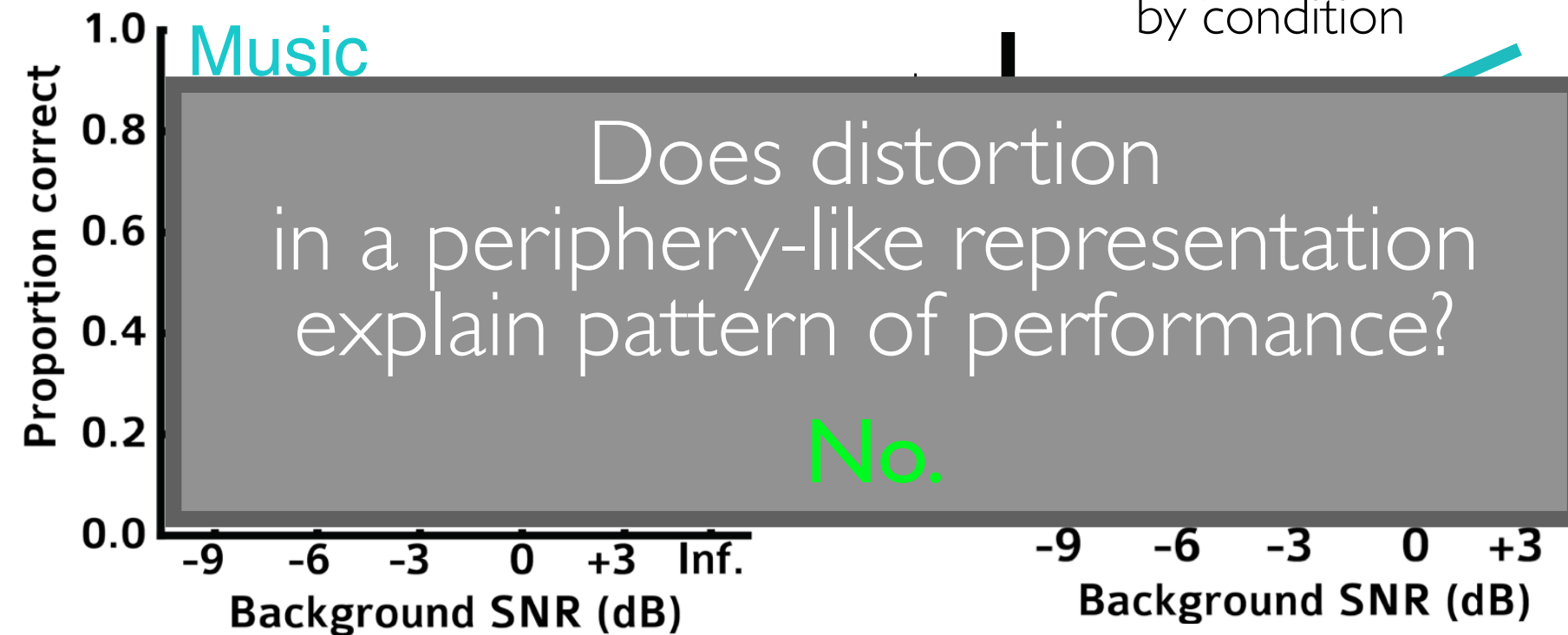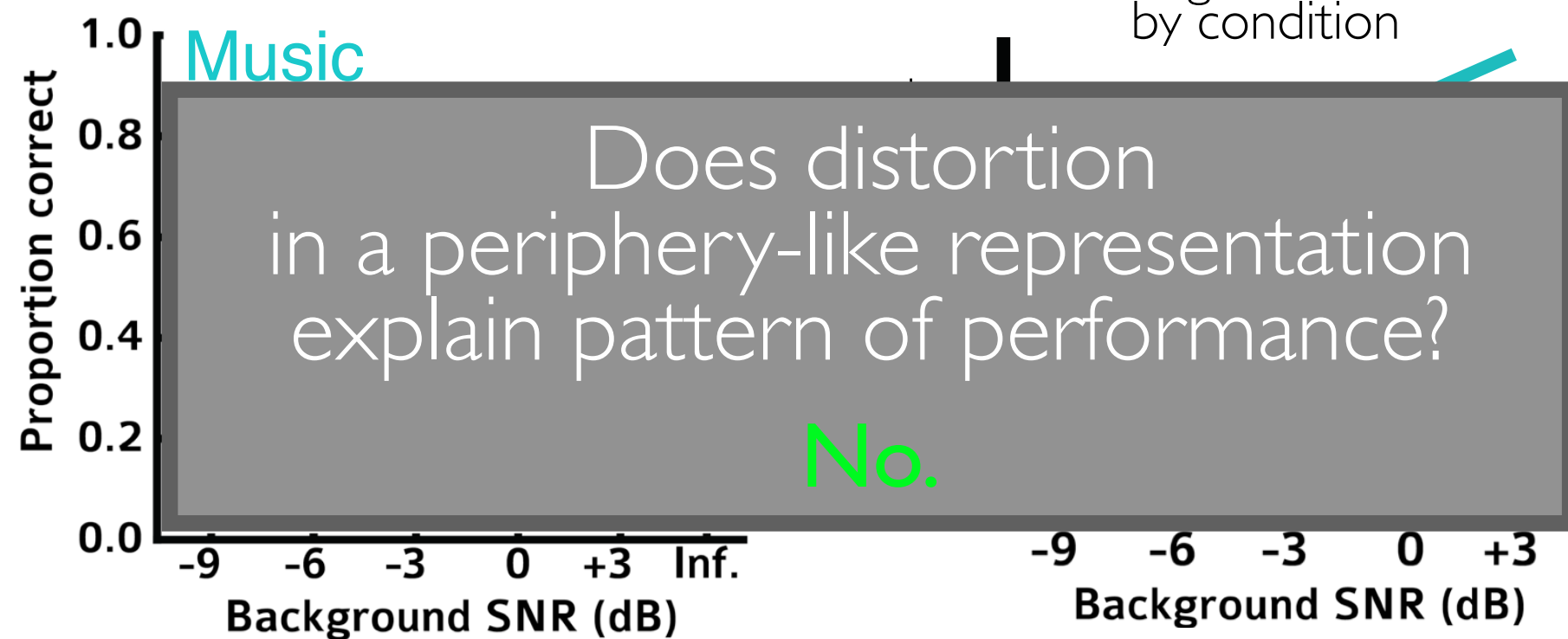
Human prop. correct

Cochleagram distortion

LEGEND

● Music   ● Auditory scenes   ● Speech babble   ● Speech-shaped noise

# Distortion and pattern of performance.



**Human Proportion Correct (n=14)**

Music

Does distortion
in a periphery-like representation
explain pattern of performance?

No.

Proportion correct

1.0
0.8
0.6
0.4
0.2
0.0

-9  -6  -3  0  +3  Inf.
Background SNR (dB)

Cochleagram distortion
by condition

-9  -6  -3  0  +3
Background SNR (dB)

Cochleagram distortion v. human
performance

Human prop. correct

Cochleagram
distortion

LEGEND
● Music    ● Auditory scenes    ● Speech babble    ● Speech-shaped noise

# Distortion and pattern of performance.



**Human Proportion Correct (n=14)**

Music

Does distortion
in a periphery-like representation
explain pattern of performance?

No.

Proportion correct

Background SNR (dB)

Cochleagram distortion
by condition

Background SNR (dB)

Cochleagram distortion v. human
performance

Human prop. correct

Cochleagram
distortion

## LEGEND

● Music    ● Auditory scenes    ● Speech babble    ● Speech-shaped noise

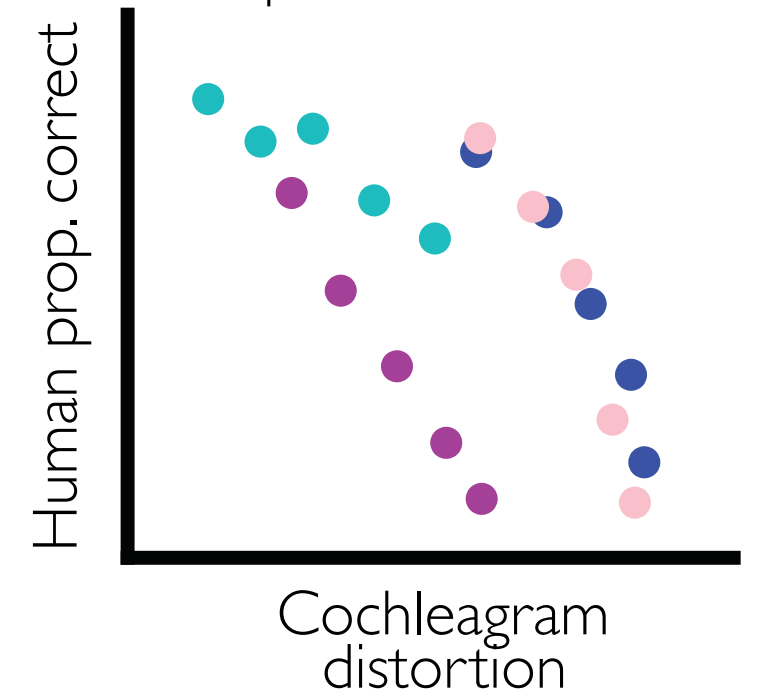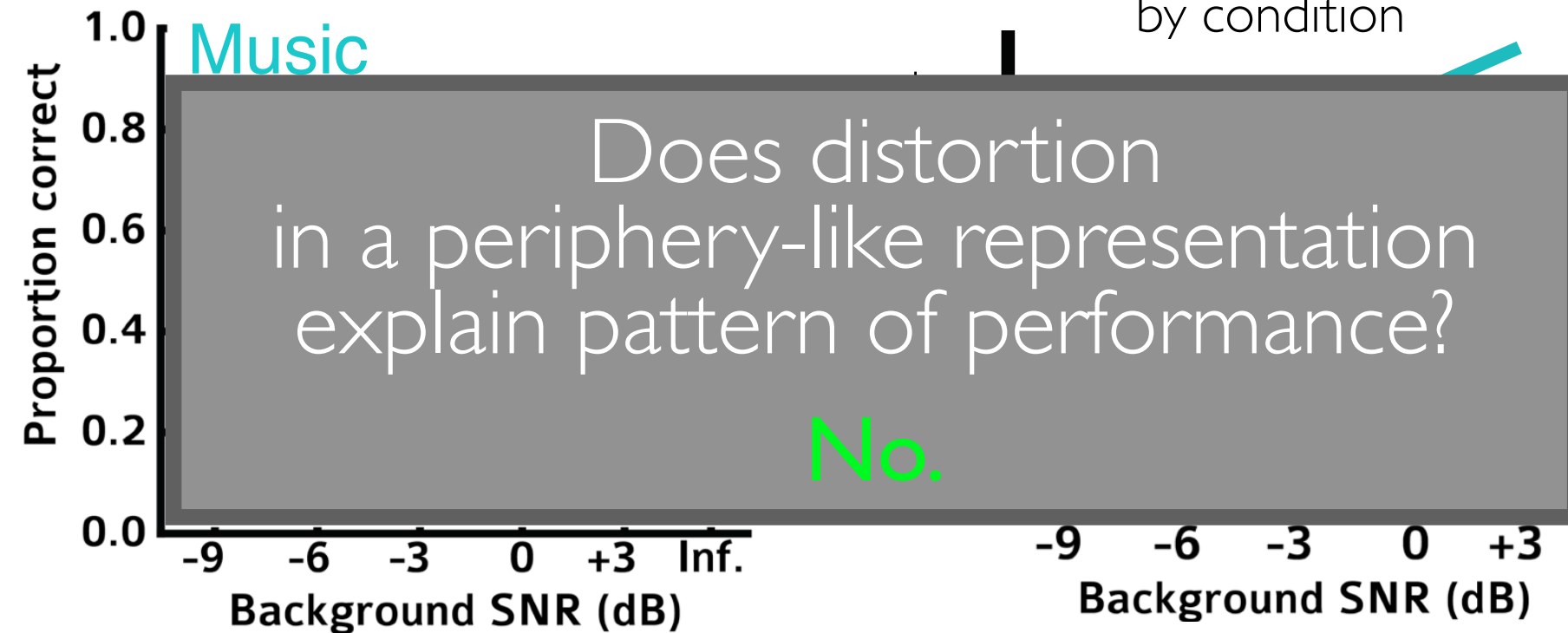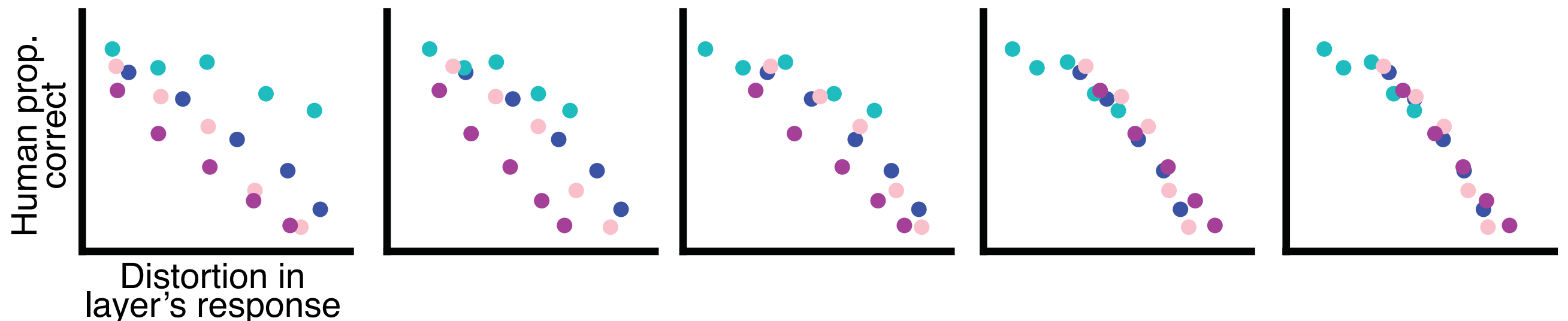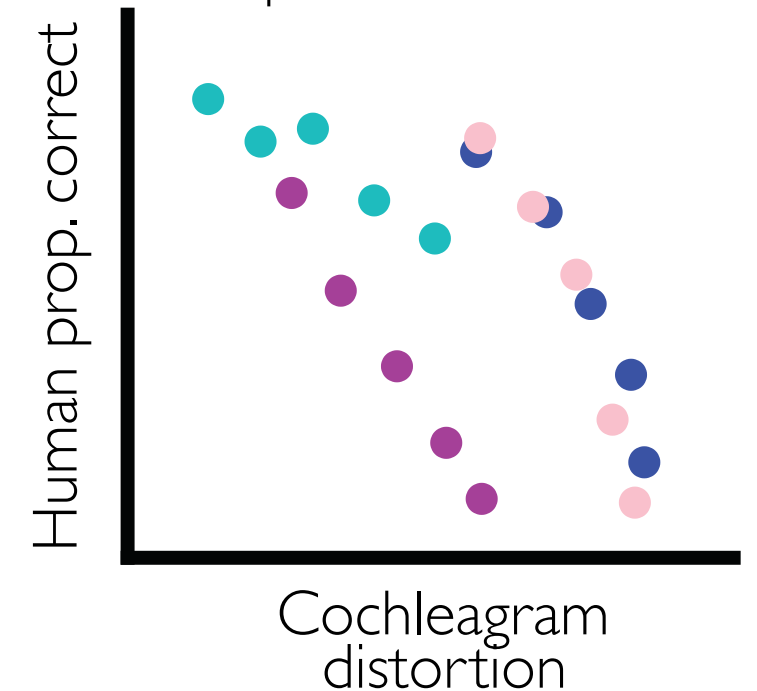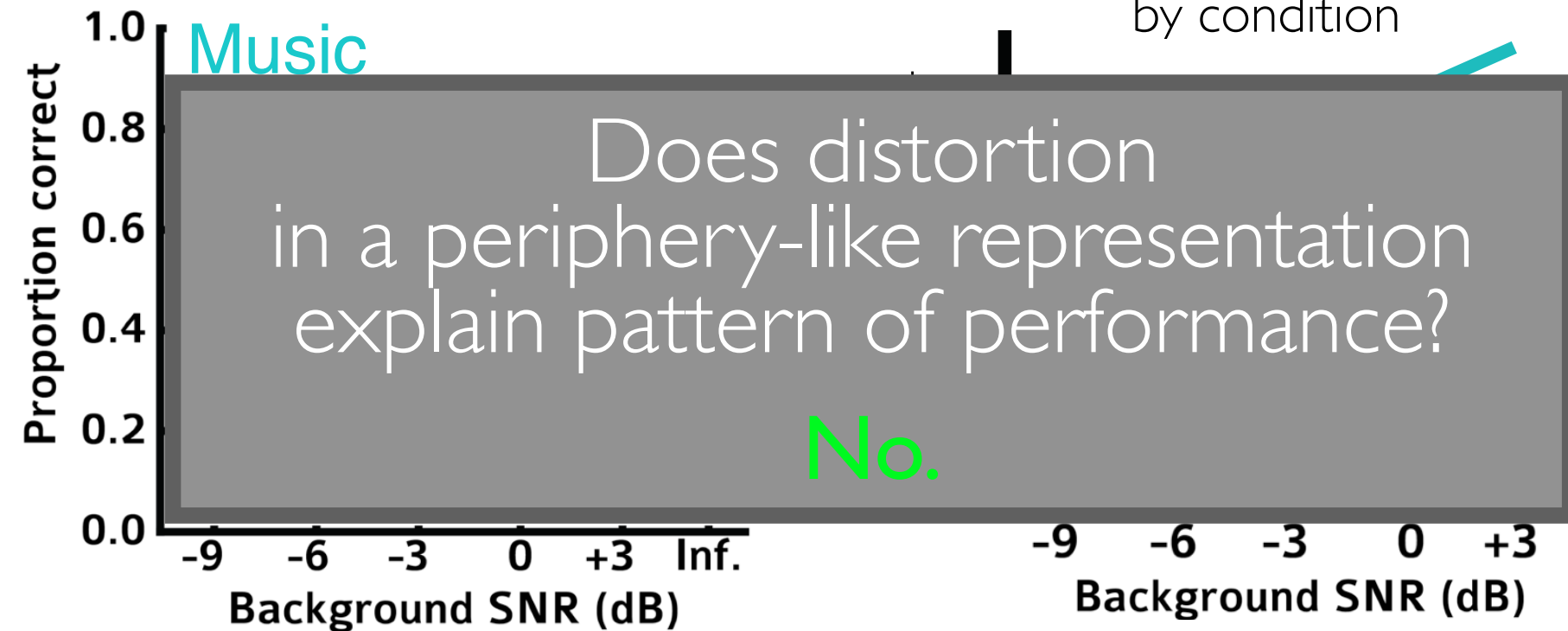Does distortion of CNN representation explain pattern of performance?

# Distortion and pattern of performance.

**Human Proportion Correct (n=14)**

Music

Does distortion
in a periphery-like representation
explain pattern of performance?

No.

Proportion correct

1.0
0.8
0.6
0.4
0.2
0.0

-9   -6   -3   0   +3   Inf.
**Background SNR (dB)**

Cochleagram distortion
by condition

-9   -6   -3   0   +3
**Background SNR (dB)**

Cochleagram distortion v. human
performance

Human prop. correct

Cochleagram
distortion

LEGEND

● Music   ● Auditory scenes   ● Speech babble   ● Speech-shaped noise

Does distortion of CNN representation explain pattern of performance?

**First layer**

Human prop.
correct

Distortion in
layer's response

# Distortion and pattern of performance.



**Human Proportion Correct (n=14)**

Music

Does distortion
in a periphery-like representation
explain pattern of performance?

No.

Proportion correct

-9  -6  -3  0  +3  Inf.
Background SNR (dB)

Cochleagram distortion
by condition

-9  -6  -3  0  +3
Background SNR (dB)

Cochleagram distortion v. human
performance

Human prop. correct

Cochleagram
distortion

LEGEND
● Music    ● Auditory scenes    ● Speech babble    ● Speech-shaped noise

Does distortion of CNN representation explain pattern of performance?

**First layer**

Human prop.
correct

Distortion in
layer's response

# Distortion and pattern of performance.



**Human Proportion Correct (n=14)**

Music

Does distortion in a periphery-like representation explain pattern of performance?

No.

Proportion correct

Background SNR (dB)
-9  -6  -3  0  +3  Inf.

Cochleagram distortion by condition

Background SNR (dB)
-9  -6  -3  0  +3

Cochleagram distortion v. human performance

Human prop. correct

Cochleagram distortion

**LEGEND**

● Music  ● Auditory scenes  ● Speech babble  ● Speech-shaped noise

## Does distortion of CNN representation explain pattern of performance?

**First layer** ⟶ **Top layer**

Human prop. correct

Distortion in layer's response

# Distortion and pattern of performance.

# Imaging Experiment

fMRI response data collected* on 165 commonly heard natural sound stimuli.

Man speaking
Flushing toilet
Pouring liquid
Tooth-brushing
Woman speaking
Car accelerating
Biting and chewing
Laughing
Typing
Car engine starting
Running water
Breathing
Keys jangling
Dishes clanking
Ringtone
Microwave
Dog barking

Road traffic
Zipper
Cellphone vibrating
Water dripping
Scratching
Car windows
Telephone ringing
Chopping food
Telephone dialing
Girl speaking
Car horn
Writing
Computer startup sound
Background speech
Songbird
Pouring water
Pop song
Water boiling

Guitar
Coughing
Crumpling paper
Siren
Splashing water
Computer speech
Alarm clock
Walking with heels
Vacuum
Wind
Boy speaking
Chair rolling
Rock song
Door knocking

•
•
•

For each voxel, measured average response to each sound:

For each voxel, measured average response to each sound:



**11065 Voxels**

**165 Sounds**

**Response Magnitude**

Data matrix:   voxels  X  sounds.

**Neural predictivity**: the ability of model to predict each individual voxel's activity using linear regression.



**11065 Voxels**

**165 Sounds**

LN

layer N?

# Model Productivity at Best Layer

Median Voxel Predictivity

~80%.

# Predictivity Difference Between High and Low Model Layers

Early layers better explanation of primary cortex, higher layers better explanation of non-primary cortex.

Significant improvement relative to existing models,

but especially in non-primary areas.

## Tonotopic
## (⊂ Primary)

# Differentiation by Region of Interest



Tonotopic
(⊂ Primary)

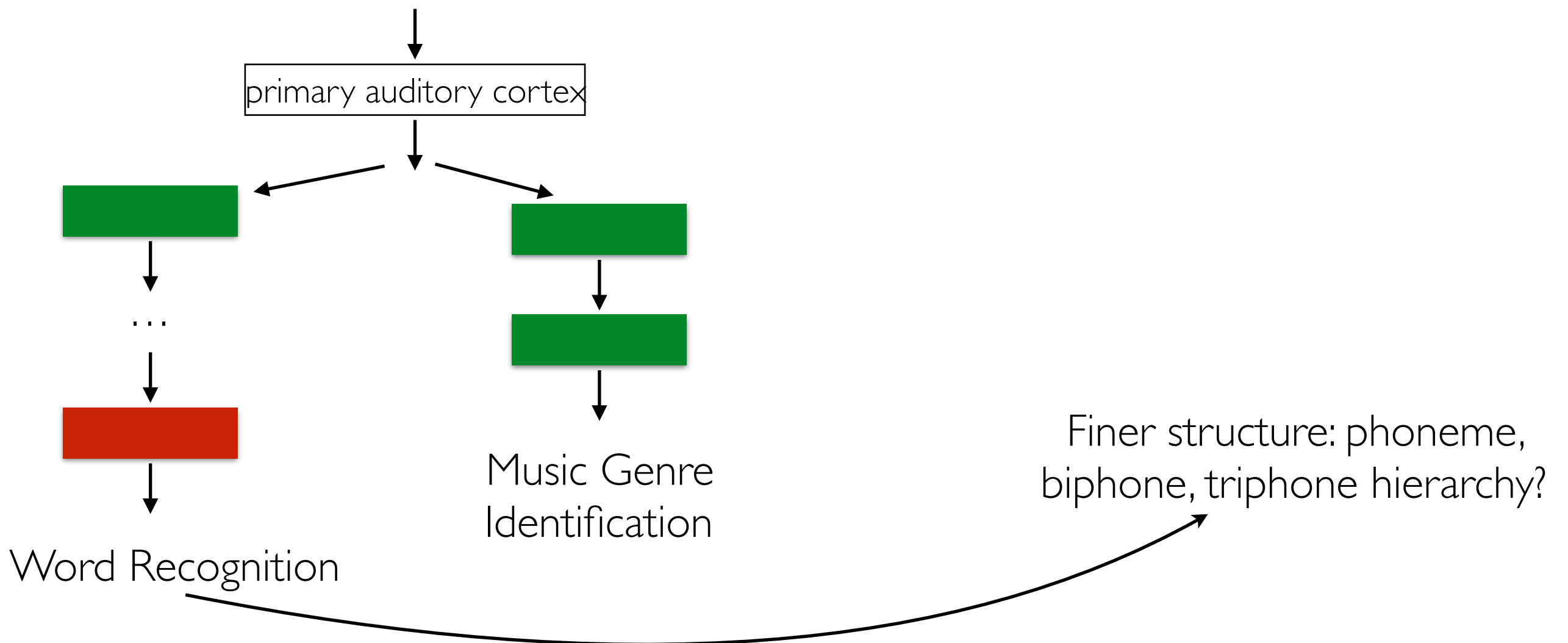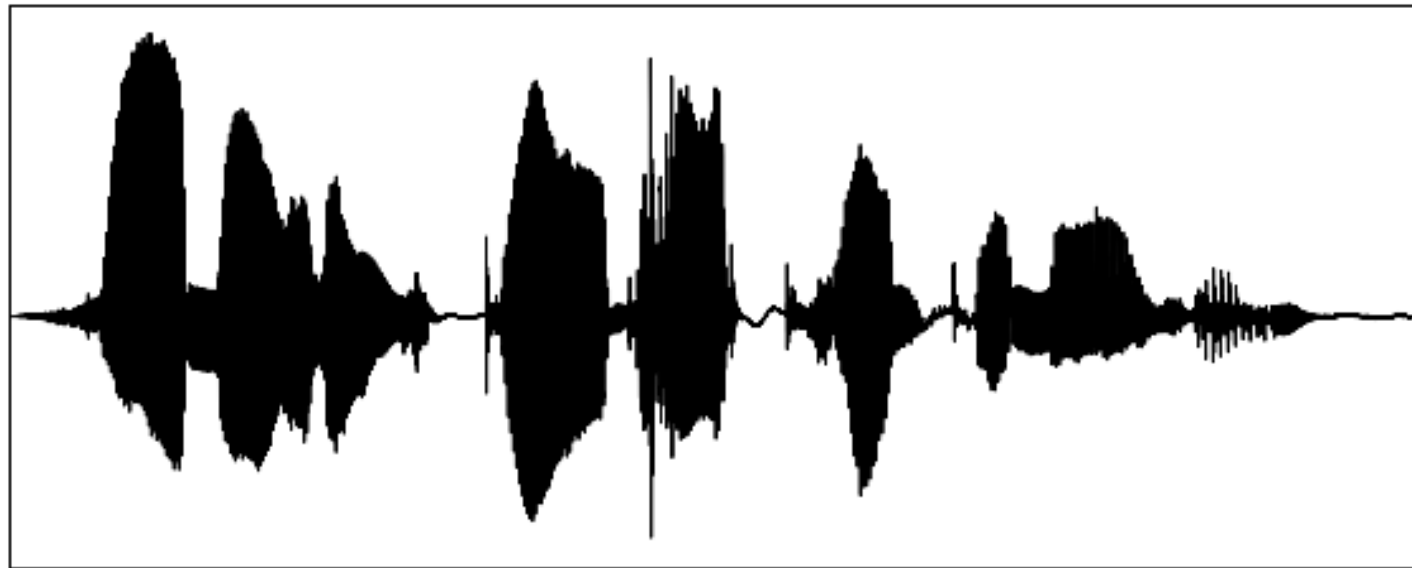Speech-selective
(⊂ Non-primary)

# Comparison of Predictivity by RoI

primary auditory cortex

Word Recognition

Music Genre
Identification

primary auditory cortex

Word Recognition

Music Genre
Identification

Finer structure: phoneme,
biphone, triphone hierarchy?

# Analysis of Model Architectures



High-variation task performance vs:

r = 0.94 ± 0.09

Auditory cortex          Higher visual cortex

Auditory cortex predictivity
*(noise-corrected voxel explained variance %)*

56.0

37.5

19.0

20          60          100

Word Recognition Performance
*(training percent correct)*

*r* = 0.87 ± 0.15

HMO

IT Explained Variance (%)

50

40

30

20

10

0

V2-like

SIFT

Different model architecture

HMAX

V1-like

PLOS09

Pixels

0.6          0.8          1.0

Categorization Performance (balanced accuracy)

Yamins et. al. (2014)

# Principle of "Goal-Driven Modeling"

visual
cortex

auditory
cortex

"Mercedes behind
Lamborghini, on a field
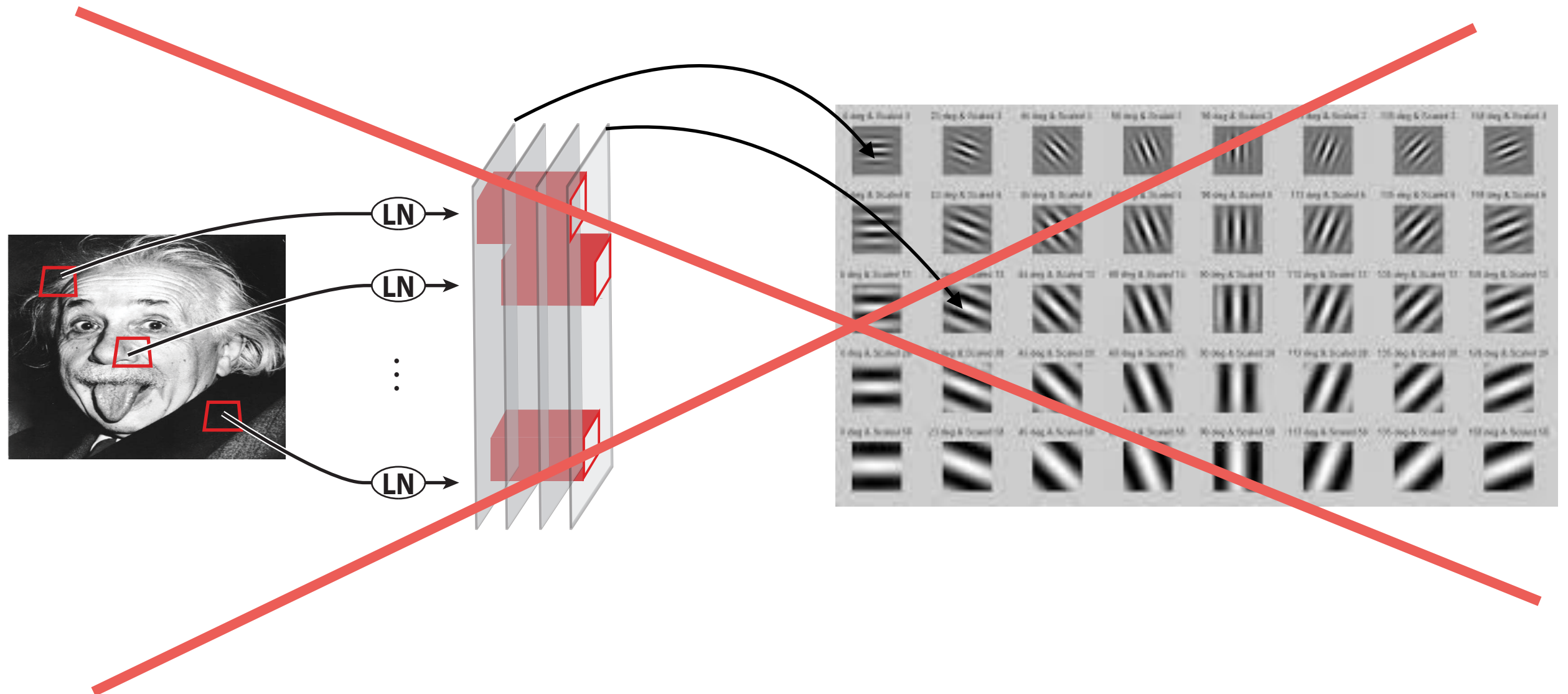in front of mountains."

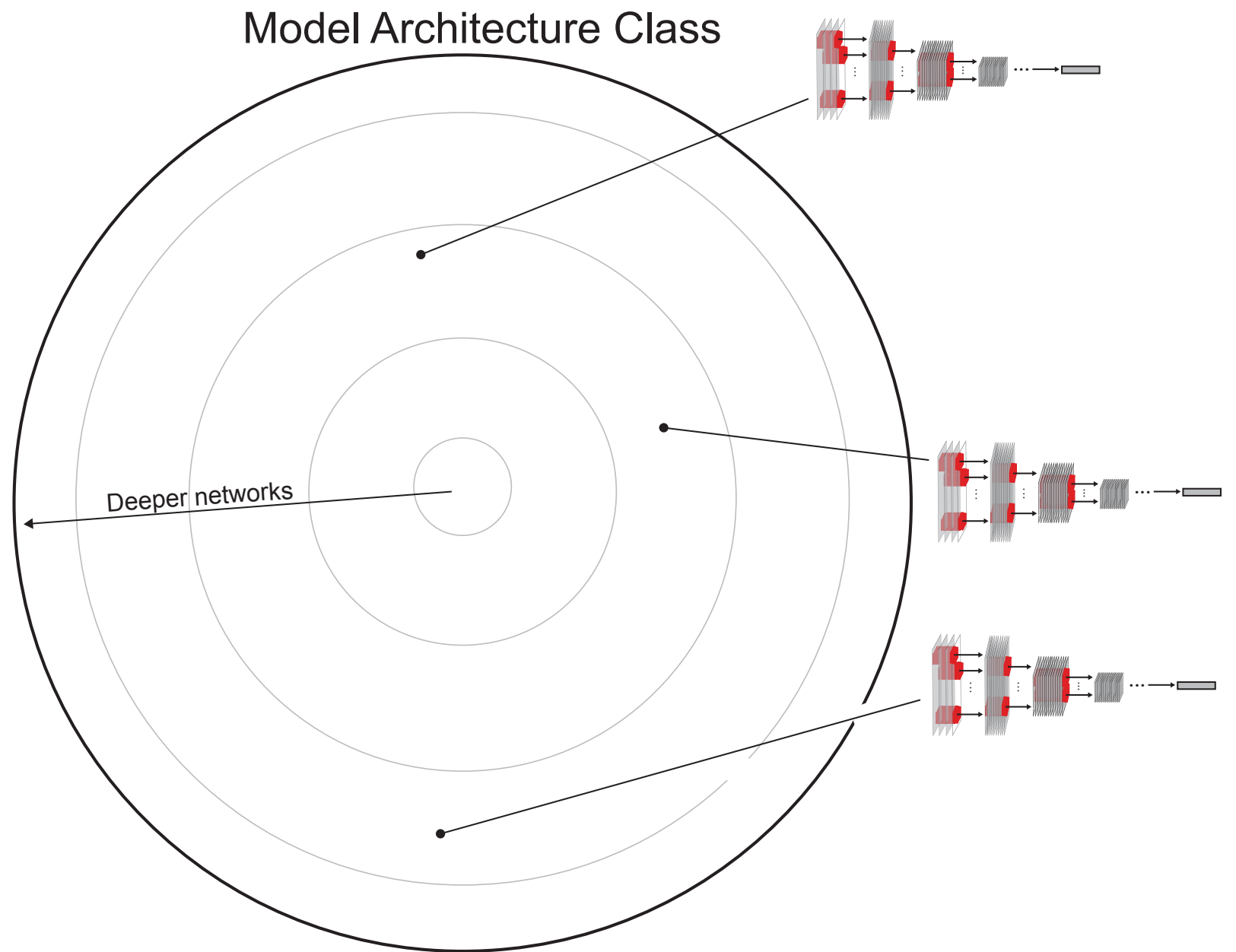"Hannah is good at
compromising"

But what type of understanding is this?

# But what type of understanding is this?



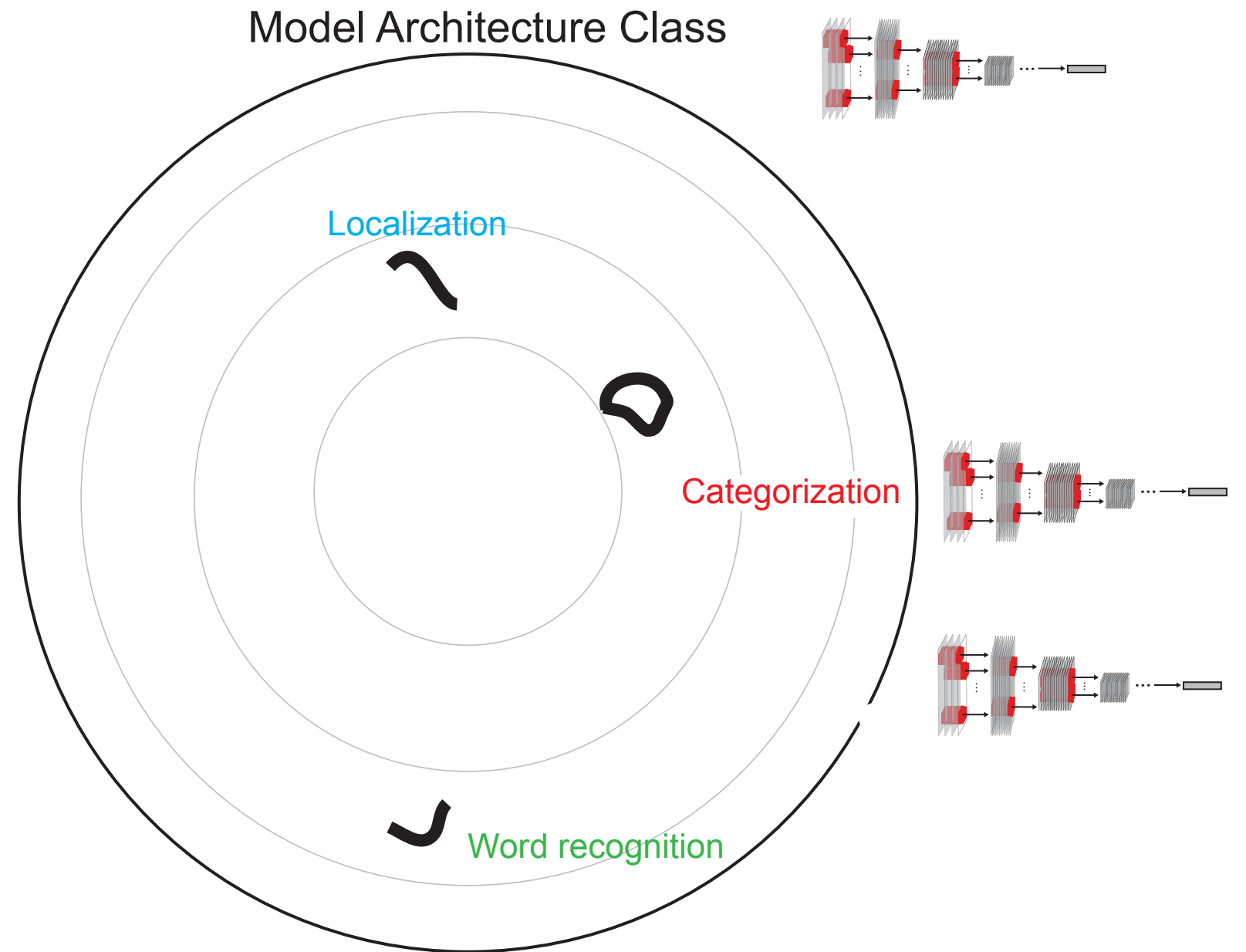*not saying this type of understanding is impossible …*

# 1. Formulate comprehensive model class (**CNNs**)
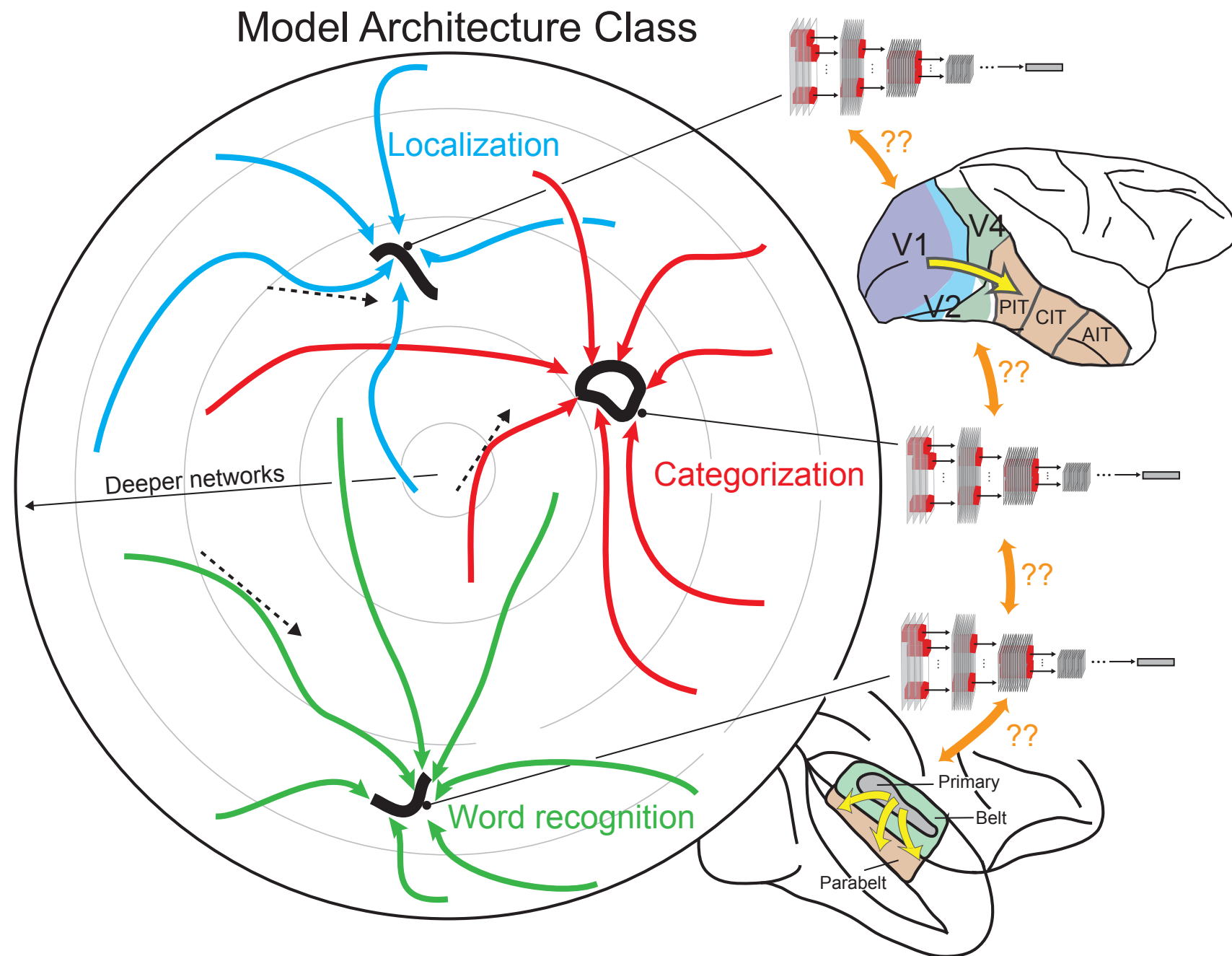


Model Architecture Class

Deeper networks

*Yamins & DiCarlo.*
***Nat. Neuro.*** *(2016)*

# 1. Formulate comprehensive model class (**CNNs**)

# 2. Choose challenging, ethologically-valid tasks (**categorization**)



Model Architecture Class

Localization

Categorization

Word recognition

*Yamins & DiCarlo.*
**Nat. Neuro.** *(2016)*

**1.** Formulate comprehensive model class (**CNNs**)

**2.** Choose challenging, ethologically-valid tasks (**categorization**)

**3.** Implement generic learning rules (**gradient descent**)



Model Architecture Class

Localization

Deeper networks

Categorization

Word recognition

*Yamins & DiCarlo.*
***Nat. Neuro.*** *(2016)*

1. Formulate comprehensive model class (**CNNs**)

2. Choose challenging, ethologically-valid tasks (**categorization**)

3. Implement generic learning rules (**gradient descent**)

> Map to brain data. (**ventral stream**)

Model Architecture Class

Localization

Categorization

Word recognition

Deeper networks

*Yamins & DiCarlo.*
**Nat. Neuro.** *(2016)*

Model Architecture Class

Localization

Deeper networks

Categorization

Word recognition

V1  V4  V2  PIT  CIT  AIT

Primary
Belt
Parabelt

??

$$\underset{a \in \mathcal{A}}{\mathrm{argmin}}[L(p_a^*)]$$

*where p\* is result of*

$$\frac{dp_a}{dt} = -\lambda(t) \cdot \langle \nabla_{p_a} L(x) \rangle_{x \in \mathcal{D}}$$

**A** = *architecture class*    **L** = *loss function*    **D** = *dataset*

Model Architecture Class

Localization

Deeper networks

Categorization

Word recognition

V1  V4  V2  PIT  CIT  AIT

Primary
Belt
Parabelt

?? ?? ?? ??

**3.**

$$\underset{a\in\mathcal{A}}{\mathrm{argmin}}[L(p_a^*)]$$

*where p\* is result of*

$$\frac{dp_a}{dt} = -\lambda(t) \cdot \langle \nabla_{p_a} L(x) \rangle_{x\in\mathcal{D}}$$

*"learning rule"*

**1.**

$\boldsymbol{A}$ = *architecture class*

**2.**

$\boldsymbol{L}$ = *loss function*      $\boldsymbol{D}$ = *dataset*

*"task"*

~~Principle~~ of "Goal-Driven Modeling"

Heuristic of "Goal-Driven Modeling"



res-net?

… after all at some point, for any given task, you'll probably "go over the hump" … perhaps when you exceed human performance or overfit on that task

dicarlolab.mit.edu ⟶ neuroailab.stanford.edu





MIT / BCS ⟶ Stanford

TEOd

TEpd

TEad

TGv
granular

*R. Lafer-Sousa and BR Conway,  Nat. Neurosci (2013)*

**Face patches**

**Color-biased regions**

Regions selective for:

- faces

- places

- bodies

- color

Where do these patches come from?

- In-born built-in structure??

- or developmentally determined by domain-specific experience?

pixel    RGC    LGN    V1    V2    V4    IT

Categorization

Categorization

pixel   RGC   LGN   V1   V2   V4   IT

Categorization

pixel　　RGC　　LGN　　V1　　V2　　V4　　IT

Categorization

pixel  RGC  LGN  V1  V2  V4  IT

Categorization

*virtual electrophysiology in models*

Shape / curvature?

"Second-order" conjunctions

pixel   RGC   LGN   V1   V2   V4   IT

Categorization

*independent goal-driven constraint?*

Shape / curvature?

"Second-order" conjunctions

pixel  RGC  LGN  V1  V2  V4  IT

Categorization

*independent goal-driven constraint?*

Shape / curvature?

"Second-order" conjunctions

pixel  RGC  LGN  V1  V2  V4  IT

Categorization

*independent goal-driven constraint?*

Shape / curvature?

"Second-order" conjunctions

pixel   RGC   LGN   V1   V2   V4   IT

*independent goal-driven constraint?*

"Second-order" conjunctions

Shape / curvature?

Categorization

pixel  RGC  LGN  V1  V2  V4  IT

Categorization

*independent goal-driven constraint?*

Shape / curvature?

"Second-order" conjunctions

Perturb:

Observe: "Behavior"

*(eg. category judgments via linear SVM)*

pixel    RGC    LGN    V1    V2    V4    IT

prediction



Perturb:

Observe: "Behavior"
*(eg. category judgments via linear SVM)*

pixel   RGC   LGN   V1   V2   V4   IT

prediction

Perturb:

Observe: "Behavior"
*(eg. category judgments via linear SVM)*

pixel  RGC  LGN  V1  V2  V4  IT

measurement

V1  V4  V2  IT

Perturb:
neural population

Observe:
actual behavior

prediction

Perturb:

Observe: units in different layer

pixel    RGC    LGN    V1    V2    V4    IT

measurement

V1   V4   V2   IT

Perturb:
neural population

Observe:
neural population

- ▸ Nonlinear? Temporal instead of rate code?

- ▸ View the process realtime?

- ▸ XX Put in Rishi learning slides

- ▸ How are they learned? What types of (e.g.) regularizations are implemented?

- ▸ Default readouts and task switching ("hyperplane management")

**H1:** Feedback implements learning of filters — form of long-term memory — not online



label?

RGC   LGN   V1   V2   V4   PIT   CIT   AIT

DOG   T(•)   ?   ?   ?

◄ - - - - -   backprop or other long-term memory error signal

**H2:** Feedback solves <u>hard cases</u> that aren't embedded in single feedforward volley …. like ambiguity.  Online inference in the ventral stream.

*(Dallenbach's cow)*



multiple inference volleys

**H3a:** Downstream, online feedback helps solve dynamic problems like (e.g.) task switching.



- ▸ category
- ▸ identity
- ▸ position
- ▸ size
- ▸ pose
- ▸ …

**H3a:** Downstream, online feedback helps solve dynamic problems like (e.g.) task switching.



stored exogenous signal

chosen

behavior

▸ category

▸ identity

▸ position

▸ size

▸ pose

▸ …

**H3a:** Downstream, online feedback helps solve dynamic problems like (e.g.) task switching.

stored exogenous signal



**chosen**

**behavior**

"Hyperplane Management"

▸ category

▸ identity

▸ position

▸ size

▸ pose

▸ …

**H3b:** Task switching algorithms reach down into ventral stream — benefit from nonlinear combinations of gating variable and IT features — e.g., attentional effects



exogenous signal

Distinguished ventrally from **H2** ("hard cases") by the nature of the task that elicits it — (e.g.) pre-cuing (volitional control) instead of (e.g.) passive viewing.

**Something about relationship to generative models**

Many parameters, **P**

How are they learned?



pixel    RGC      LGN      V1      V2      V4      IT

Gradient descent eq:

$$\frac{dp}{dt} = -\lambda(t) \cdot \frac{\partial L}{\partial P}$$

L = loss function

$\lambda$ = learning rate

Many parameters, **P**

How are they learned?



pixel    RGC         LGN         V1         V2         V4         IT

Gradient descent eq:

$$\frac{dp}{dt} = -\lambda(t) \cdot \frac{\partial L}{\partial P}$$

L = loss function

$\lambda$ = learning rate

In current standard practice:

L = soft-max loss computed relative large numbers of externally-provided semantic labels.

Many parameters, **P**

How are they learned?



pixel   RGC      LGN      V1      V2      V4      IT

Gradient descent eq:

$$\frac{dp}{dt} = -\lambda(t) \cdot \frac{\partial L}{\partial P}$$

L = loss function

$\lambda$ = learning rate

In current standard practice:

L = loss computed via large numbers of externally-provided semantic labels.

Ideally:

L = un- or semi-supervised function computable from easily accessible data about agent's environment

Many parameters, **P**

How are they learned?



pixel    RGC    LGN    V1    V2    V4    IT

Gradient descent eq:

$$\frac{dp}{dt} = -\lambda(t) \cdot \frac{\partial L}{\partial P}$$

L = loss function

$\lambda$ = learning rate

(1) Which parameters are learned vs developed or evolved?

(2) What are the right loss functions(s)?

(3) How are the loss functions and the GDE implemented/ approximated via neural circuits?

Can goal-drive modeling approaches generalize to other areas?

For example, in auditory cortex:

‣ can HCNN models explain higher auditory cortex?

‣ which tasks best explain functional organization of AC?

‣ how to auditory-optimized architectures related to visual ones?

# Core / Belt / Parabelt Structure



Core area

Belt area

A1

P

A

Parabelt area

*monkey
*

*Tramo et. al, Curr. Opin. Neuro. (1999)*

Spatiotemporal filtering?

Core area

Belt area

???



A1

P

A

*Tramo et. al, Curr. Opin. Neuro. (1999)*

\*monkey

\*

Parabelt area

???

Example: use computational models to help deepen understanding of non-primary areas.

**Task-Driven Modeling:**

1. Optimize for performance on a challenging auditory task (600-way work recognition in noisy speech)

2. Compare to neural data.

Apply to auditory tasks, where the regions themselves are less well known.

# Which layer best predicts each voxel's responses?



CNN suggests hierarchical functional organization of auditory cortex.

Higher layer

Layer 6
Layer 5
Layer 4
Layer 3
Layer 2

Lower layer

**Primary** auditory cortex: predicted by **lower CNN layers**.
**Non-primary** auditory cortex: predicted by **higher CNN layers**.

# Analysis of Model Architectures



High-variation task performance vs:

r = 0.94 ± 0.09

Auditory cortex

Higher visual cortex

Auditory cortex predictivity
*(noise-corrected voxel explained variance %)*

56.0

37.5

19.0

20    60    100

Word Recognition Performance
*(training percent correct)*

*r* = 0.87 ± 0.15

IT Explained Variance (%)

50

40

30

20

10

0

HMO

V2-like

SIFT

Different model architecture

HMAX

PLOS09

V1-like

Pixels

0.6    0.8    1.0

Categorization Performance (balanced accuracy)

Yamins et. al. (2014)

Many visual behaviors beyond vision at a glance (e.g.):

- ▸ Scene understanding over multiple saccades

- ▸ Strategic decision-making in complex environments

Involve integration of working memory, likely via RNNs.

Match between models and data at category confusion level is pretty good …

Match between models and data at category confusion level is pretty good …



work of:

Elias Issa     Rishi Rajalingham     Kohitij Kar     Kailyn Schmidt     Jim DiCarlo

# Digging deeper into understanding visual cortex

Match between models and data at category confusion level is pretty good … but less good at *image* grain:



work of:

Elias Issa

Rishi Rajalingham

Kohitij Kar

Kailyn Schmidt

Jim DiCarlo

# Digging deeper into understanding visual cortex

Match between models and data at category confusion level is pretty good … but less good at *image* grain:



**Object grain** — x-axis: Performance (relative to human); y-axis: Object-level Consistency. Labels: Monkey, GoogLeNet (v3) synthetic trained, HMAX, V1.

**Image grain** — x-axis: Performance (relative to human); y-axis: g. Labels: Monkey, GoogLeNet (v3) synthetic trained, HMAX, V1.

*(remember, neural fits only ~50%)…*

work of:

Elias Issa     Rishi Rajalingham     Kohitij Kar     Kailyn Schmidt     Jim DiCarlo

**3.**

$$\underset{a \in \mathcal{A}}{\mathrm{argmin}}[L(p_a^*)]$$

*where p\* is result of*

$$\frac{dp_a}{dt} = -\lambda(t) \cdot \langle \nabla_{p_a} L(x) \rangle_{x \in \mathcal{D}}$$

*"learning rule"*

**1.**

**A** = *architecture class*

**2.**

**L** = *loss function*    **D** = *dataset*

*"task"*

Model Architecture Class

Localization

Deeper networks

Categorization

Word recognition

V1  V4  V2  PIT  CIT  AIT

Primary  Belt  Parabelt

?? ?? ?? ??

**3.**

$$\operatorname*{argmin}_{a \in \mathcal{A}}[L(p_a^*)]$$

*where p\* is result of*

$$\frac{dp_a}{dt} = -\lambda(t) \cdot \langle \nabla_{p_a} L(x) \rangle_{x \in \mathcal{D}}$$

*"learning rule"*

**1.**

**A** = *architecture class*

**2.**

**L** = *loss function*     **D** = *dataset*

*"task"*

Three hypotheses:

1) the task (loss function **L** or dataset **D**) is wrong

2) the architecture class (**A**) is wrong

3) ~~the learning rule (argmin, SGD, &c) rule is wrong~~
*bet: some version of approximate backprop-like error correction is reasonable*

Optimize models of the current structure to directly
match the neural data …

Optimize models of the current structure to directly match the neural data …

no                    yes

Model class is
wrong

Task
is wrong

Optimize models of the current structure to directly match the neural data …

no          yes

Model class is wrong                    Task is wrong

… but not enough neural data?

Optimize models of the current structure to directly match the behavioral data … then check against neural data.

no

yes

Model class is wrong

Task is wrong

Optimize models of the current structure to directly match the behavioral data … then check against neural data.

no

yes

Model class is wrong

Task is wrong

Eli Wang

Elias Issa

Rishi Rajalingham

Kohitij Kar

Kailyn Schmidt

Jim DiCarlo

Optimize models of the current structure to directly match the behavioral data … then check against neural data.

(i) predict vector of errors



imagenet

softmax

regression

error pattern

layer N?

Eli
Wang

Elias
Issa

Rishi
Rajalingham

Kohitij
Kar

Kailyn
Schmidt

Jim
DiCarlo

Optimize models of the current structure to directly match the behavioral data … then check against neural data.

(ii) as actual error pattern



Eli
Wang

Elias
Issa

Rishi
Rajalingham

Kohitij
Kar

Kailyn
Schmidt

Jim
DiCarlo

Optimize models of the current structure to directly match the behavioral data … then check against neural data.

(iii) as multiplier — indicator of niche?



Eli
Wang

Elias
Issa

Rishi
Rajalingham

Kohitij
Kar

Kailyn
Schmidt

Jim
DiCarlo

# Better tasks (loss functions)

*less normative "task"*                                         *more normative "task"*

$\longleftrightarrow$

Fit neural data

# Better tasks (loss functions)

*less normative "task"*                                        *more normative "task"*

⟵──────────────────────────────⟶

Fit neural data

Fit categorization
error pattern
…
check against neural
data

# Better tasks (loss functions)

*less normative "task"*                                                        *more normative "task"*

$\longleftrightarrow$

Fit neural data

Fit categorization
error pattern
…
check against neural
data

Solve non-
categorization
tasks
…
check against
neural data

# Better tasks (loss functions)

*less normative "task"* $\longleftrightarrow$ *more normative "task"*

Fit neural data

Fit categorization error pattern

…

check against neural data

Solve non-categorization tasks

…

check against neural data

But which non-categorical tasks?

# Better tasks (loss functions)

Pose / position estimation

Normal/Depth estimation

Segmentation

Chengxu Zhuang

Future prediction under agent-controlled actions

$x_0$

$x_1$

$x_2$

$\hat{x}_0$

$\hat{x}_1$

$\hat{x}_2$

working memory network

Nick Haber

Damian Mrowca

Fei-Fei Li

Where should the tasks be imposed? (intermediate?)



*pose =
(0, 45, 0)*

*approximate
surface
normal map*

category = "bald eagle"

*shallow category-dedicated network*

Strategy: optimize over architectures for solving
joint tasks, compare to neural data

*How much work can less heavily supervised tasks do?*

# Comparing to Neural Data

✳ Various second-order metrics: encoding regression, RSAs, etc

✳ Behavioral consistency — pattern of errors at various grains of detail

✳ But really, there is a developmental hypothesis implicit in these models. Time course of all metrics should be matched:



Layer-1 Filters at t = 0

... at t = 1

... at t = 40

Cateorizaiton Performance

0.75

0.55

0.35

Training Timecourse (*thousands of iterations*)

*Use developmental data separate more biologically correct loss functions from less correct ones?*

# Better architectures



AlexNet pool5 (1000 components) neural fits

work of

Kohitij
Kar

Jonas
Kubilius

Jim
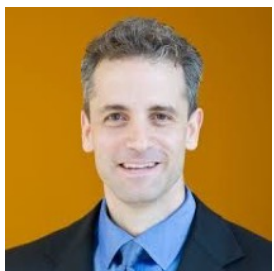DiCarlo

**Aran Nayebi**

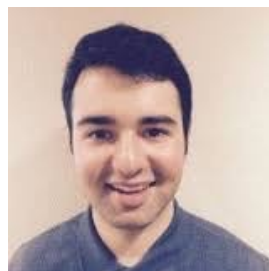Surya Ganguli

**Jonas Kubilius**

Maryann Rui

**Kohitij Kar**

Jim DiCarlo

# Better architectures



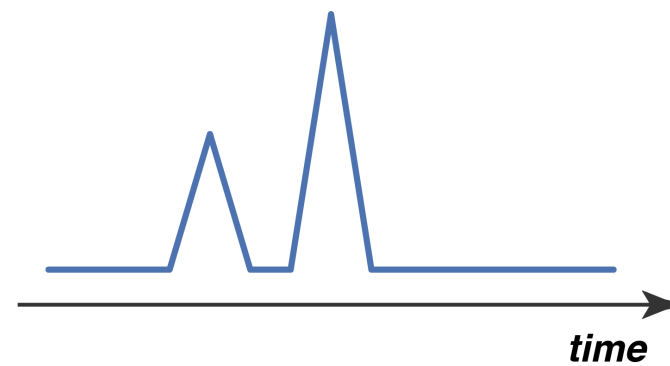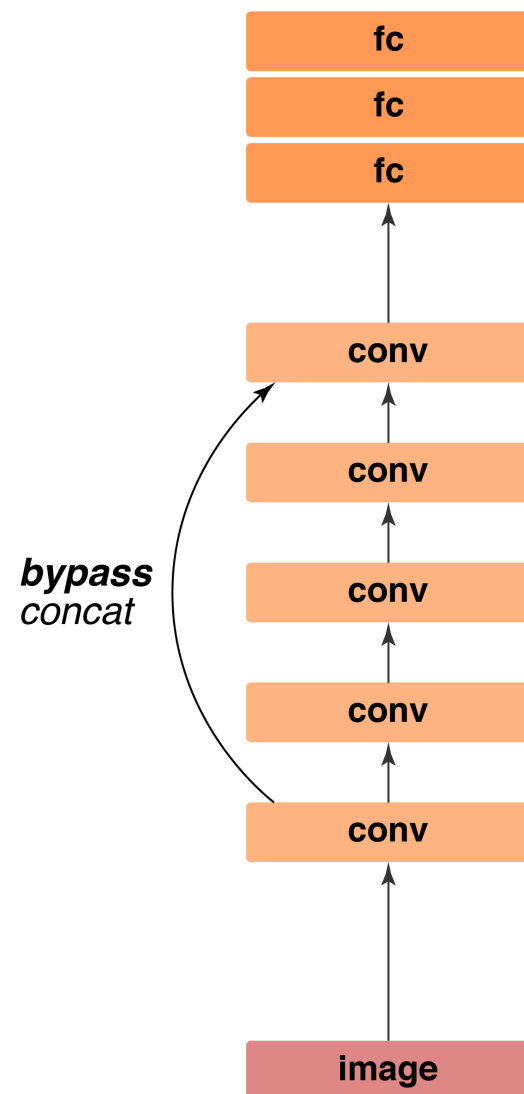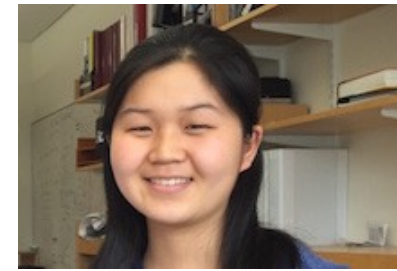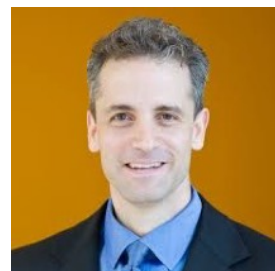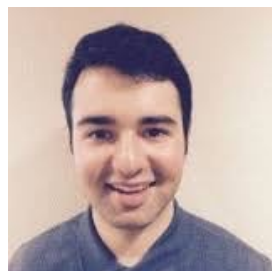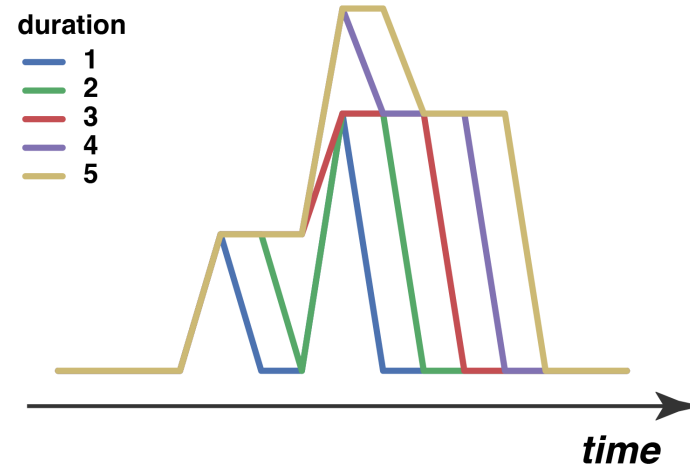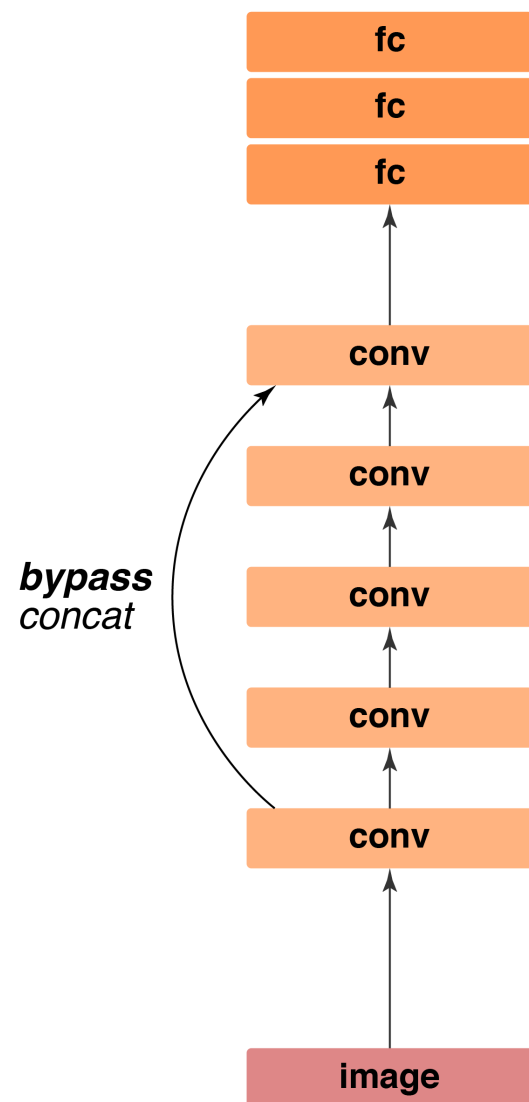Aran **Nayebi**  Surya Ganguli  **Jonas Kubilius**  Maryann Rui  **Kohitij Kar**  Jim DiCarlo
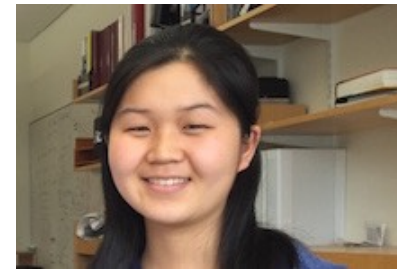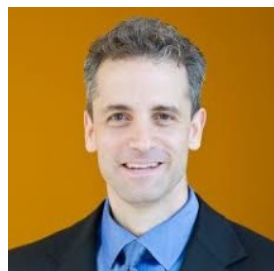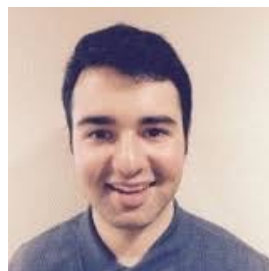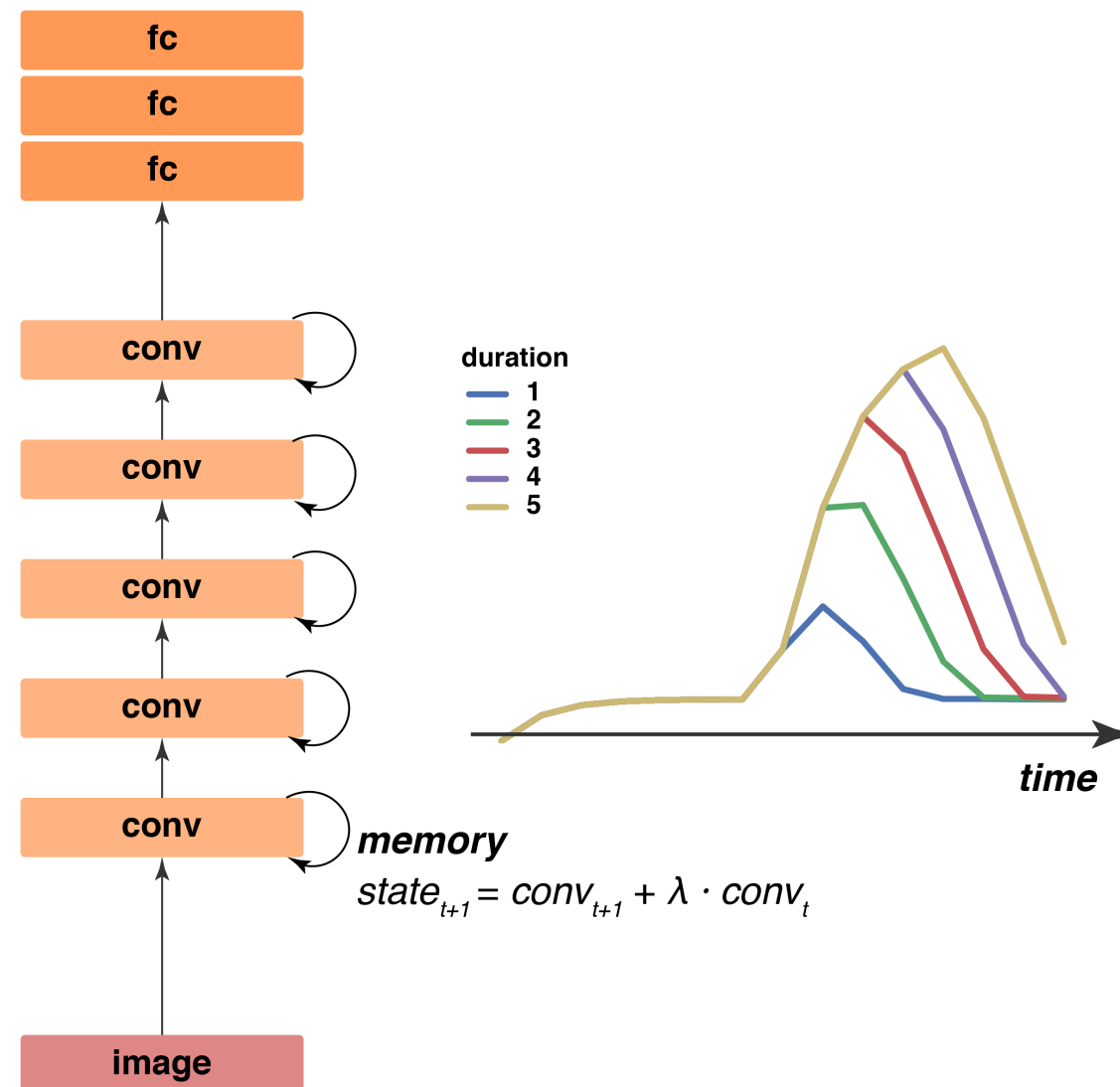
**Aran Nayebi**    Surya Ganguli    **Jonas Kubilius**    Maryann Rui    **Kohitij Kar**    Jim DiCarlo

# Better architectures

# Better architectures



bypass
concat

fc
fc
fc
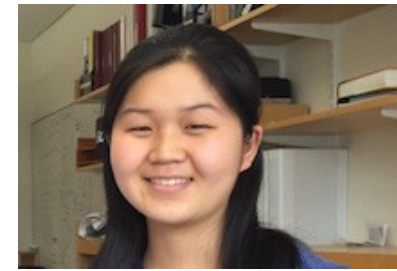conv
conv
conv
conv
conv
image

**Aran Nayebi**  
Surya Ganguli  
**Jonas Kubilius**  
Maryann Rui  
**Kohitij Kar**  
Jim DiCarlo

# Better architectures



**bypass** *concat*

fc
fc
fc
conv
conv
conv
conv
conv
image

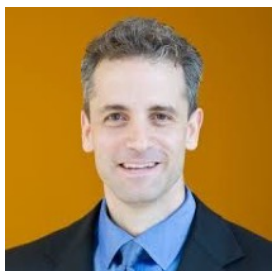**Aran Nayebi**  Surya Ganguli  **Jonas Kubilius**  Maryann Rui  **Kohitij Kar**  Jim DiCarlo
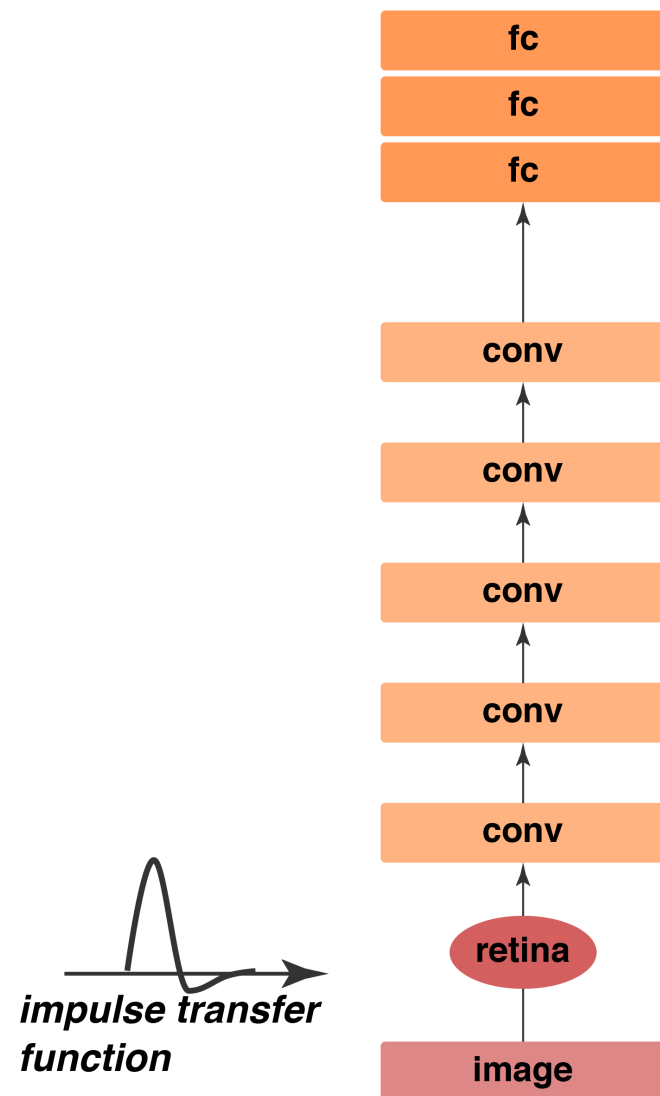
# Better architectures



fc

fc

fc

conv

conv

**bypass**
*concat*

conv

conv

conv

image

**Aran Nayebi**

Surya Ganguli

**Jonas Kubilius**

Maryann Rui

**Kohitij Kar**

Jim DiCarlo

# Better architectures



fc

fc

fc

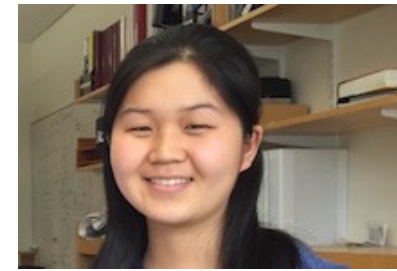conv

*bypass*
concat

conv

conv

conv

conv

image

**Aran
Nayebi**

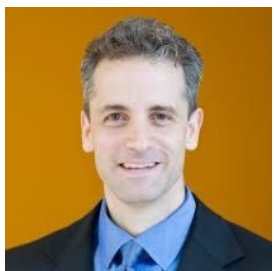Surya
Ganguli

**Jonas
Kubilius**

Maryann
Rui

**Kohitij
Kar**

Jim
DiCarlo

# Better architectures



fc

fc

fc

conv

conv

**bypass**
*concat*

conv

conv

conv

image

**Aran Nayebi**

Surya Ganguli

**Jonas Kubilius**

Maryann Rui

**Kohitij Kar**

Jim DiCarlo

# Better architectures

# Better architectures



bypass
concat

Aran
Nayebi

Surya
Ganguli

Jonas
Kubilius

Maryann
Rui

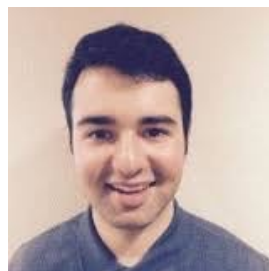Kohitij
Kar

Jim
DiCarlo

# Better architectures



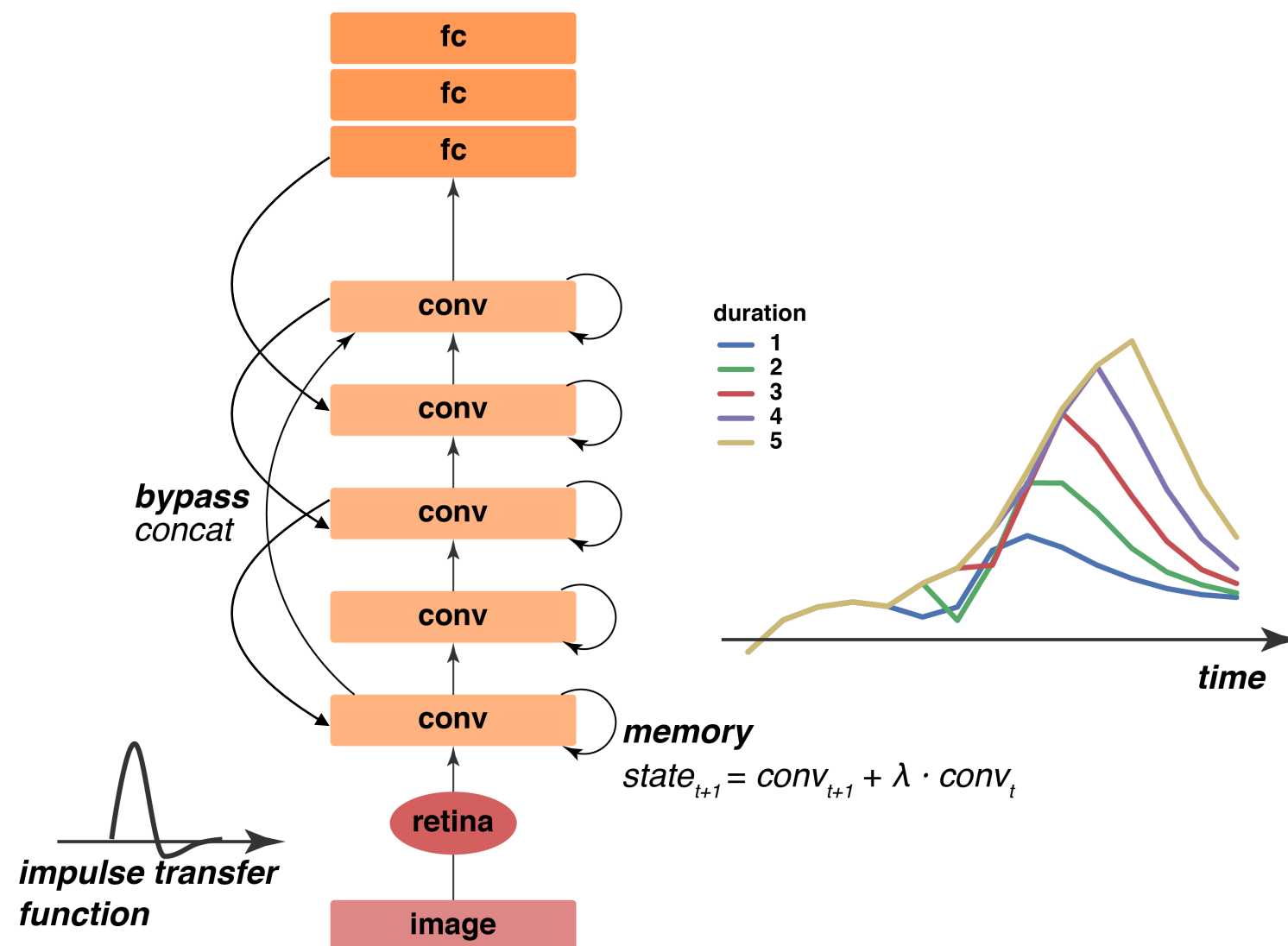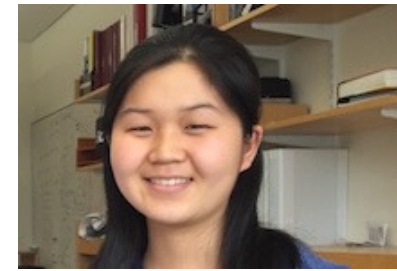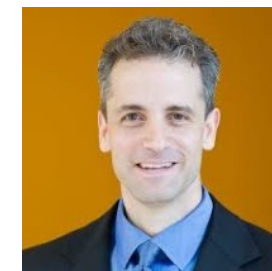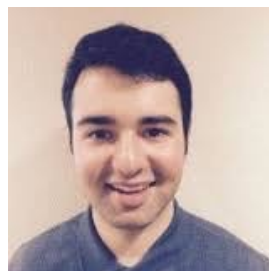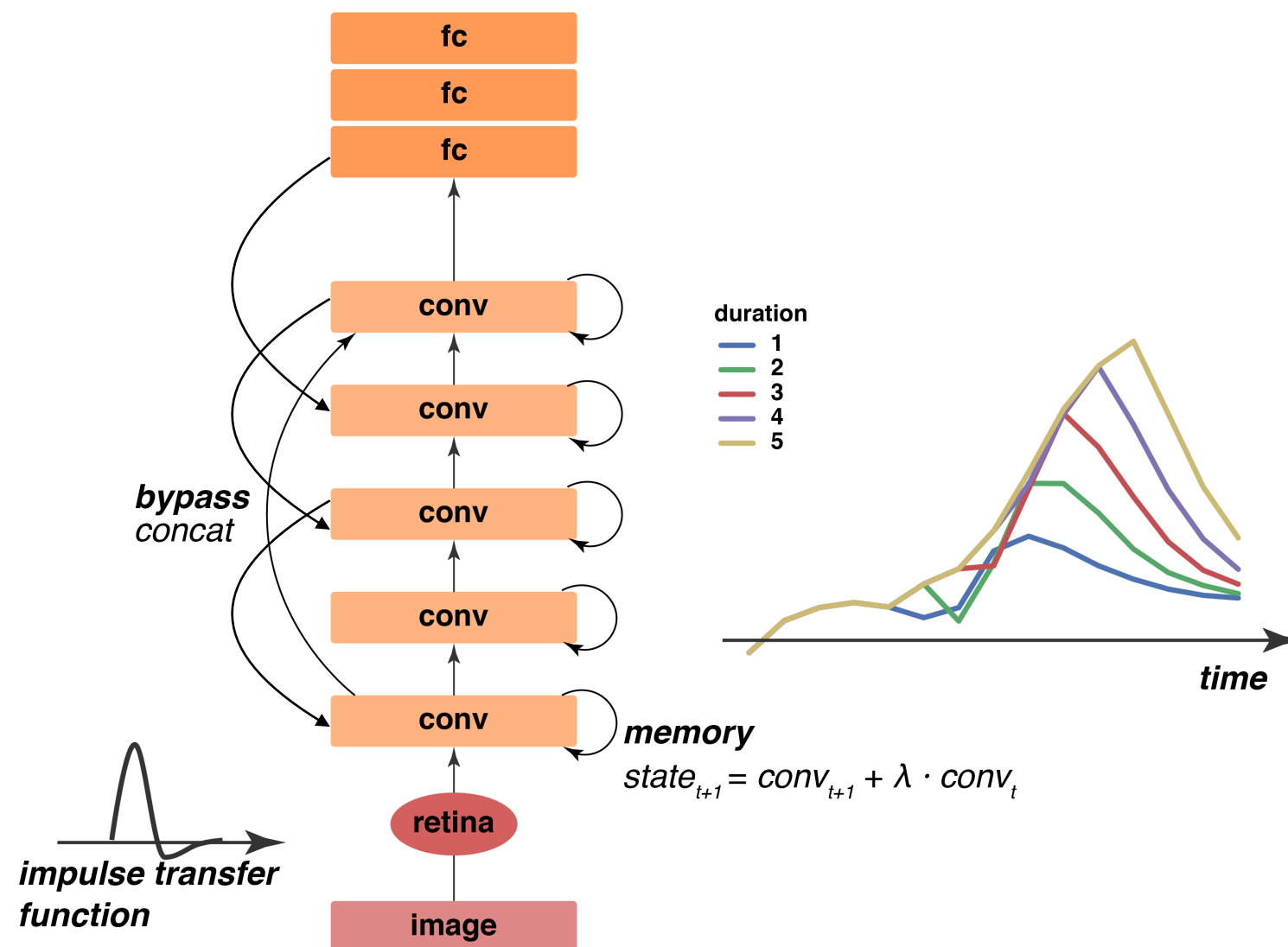Aran **Nayebi**  Surya Ganguli  Jonas **Kubilius**  Maryann Rui  Kohitij **Kar**  Jim DiCarlo

fc

fc

fc

conv

conv

*bypass*
concat

conv

conv
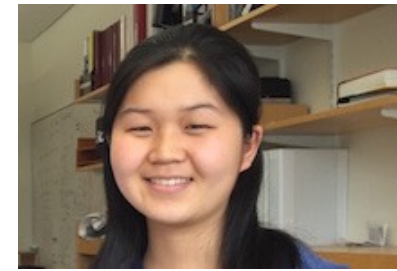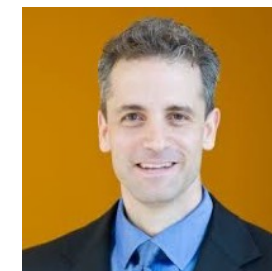
conv

image

duration
1
2
3
4
5

time

**Aran
Nayebi**

Surya
Ganguli

**Jonas
Kubilius**

Maryann
Rui

**Kohitij
Kar**

Jim
DiCarlo

# Better architectures



fc
fc
fc

conv
conv
conv
conv
conv

*memory*

$$state_{t+1} = conv_{t+1} + \lambda \cdot conv_t$$

image

duration
1
2
3
4
5

*time*

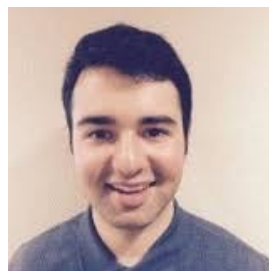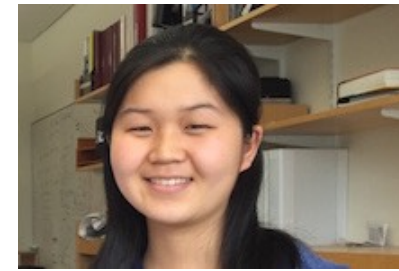**Aran Nayebi**  Surya Ganguli  **Jonas Kubilius**  Maryann Rui  **Kohitij Kar**  Jim DiCarlo

# Better architectures

# Better architectures



impulse transfer function

duration
1
2
3
4
5

time

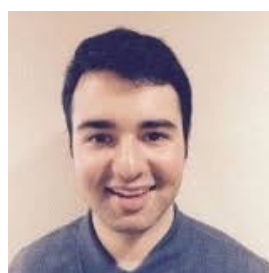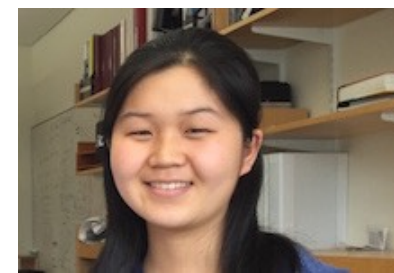bypass concat

memory

$$state_{t+1} = conv_{t+1} + \lambda \cdot conv_t$$

**Aran Nayebi**   Surya Ganguli   **Jonas Kubilius**   Maryann Rui   **Kohitij Kar**   Jim DiCarlo

**fc**

**fc**

**fc**

**conv**

**conv**

**conv**

**conv**

**conv**

*bypass*
concat

*memory*
$state_{t+1} = conv_{t+1} + \lambda \cdot conv_t$

duration
1
2
3
4
5

*time*

**retina**

*impulse transfer
function*

**image**

**Aran
Nayebi**

Surya
Ganguli

**Jonas
Kubilius**

Maryann
Rui

**Kohitij
Kar**

Jim
DiCarlo

fc
fc
fc

conv

conv

**bypass**
*concat*

conv

conv

conv

**memory**
$state_{t+1} = conv_{t+1} + \lambda \cdot conv_t$

retina

**impulse transfer
function**

image

**duration**
— 1
— 2
— 3
— 4
— 5

**time**

What task(s)?

a) vanilla categorization

**Aran
Nayebi**

Surya
Ganguli

**Jonas
Kubilius**

Maryann
Rui

**Kohitij
Kar**

Jim
DiCarlo

# Better architectures



fc

fc

fc

conv

conv

**bypass**
*concat*

conv

conv

conv **memory**

$state_{t+1} = conv_{t+1} + \lambda \cdot conv_t$

**retina**

*impulse transfer function*
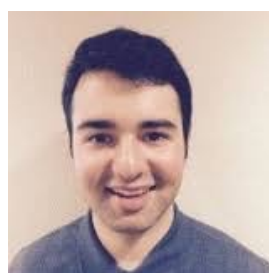
**image**

duration
— 1
— 2
— 3
— 4
— 5

**time**

What task(s)?

a) vanilla categorization

b) time-discounting

$$L = \sum_t \gamma^t \cdot L_t$$

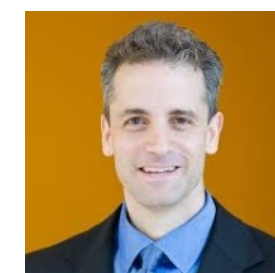_be accurate but also fast_

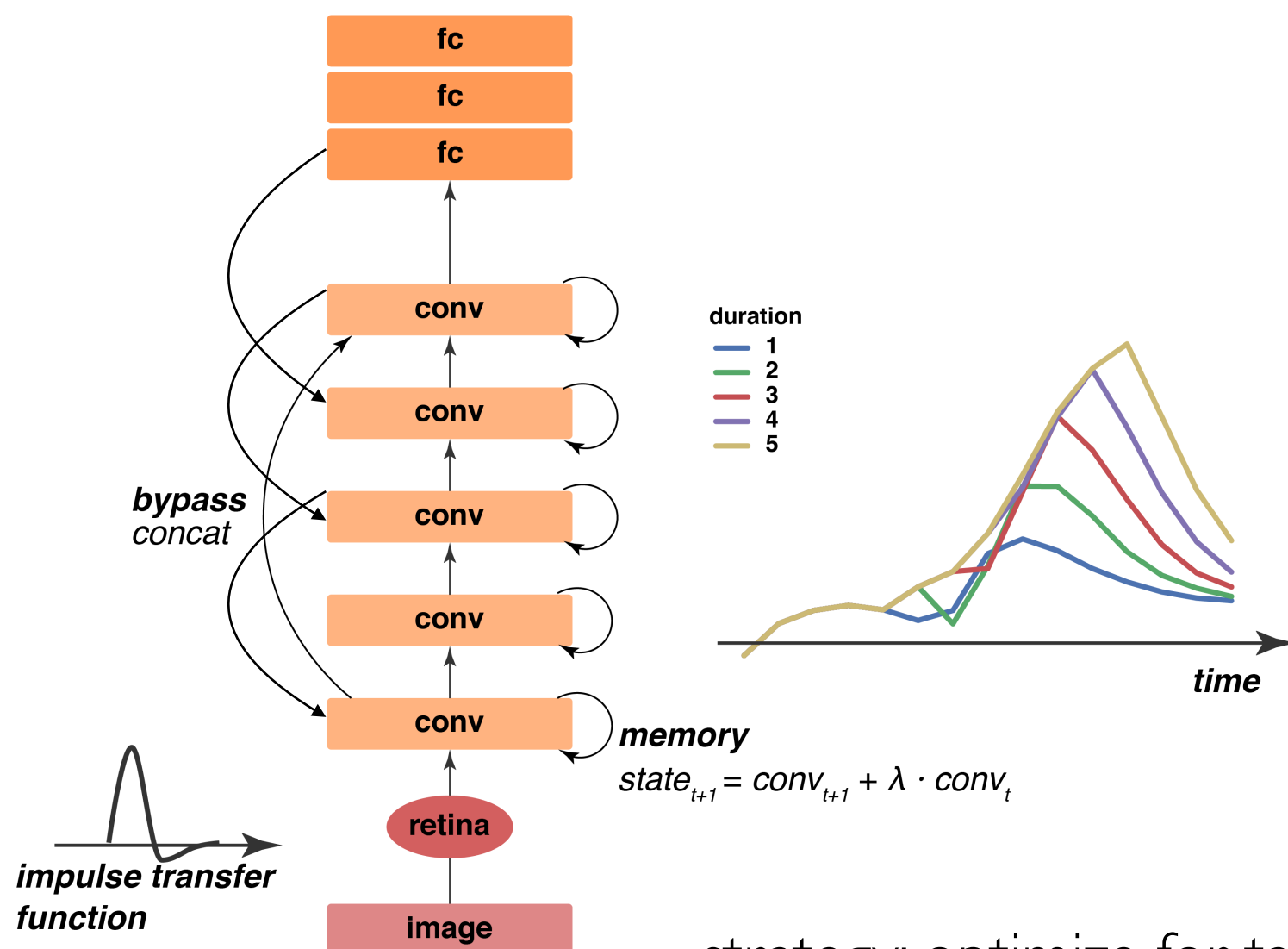**Aran Nayebi**    Surya Ganguli    **Jonas Kubilius**    Maryann Rui    **Kohitij Kar**    Jim DiCarlo

# Better architectures

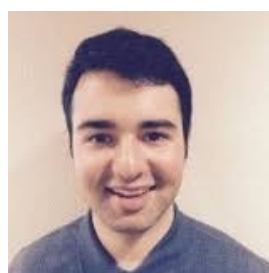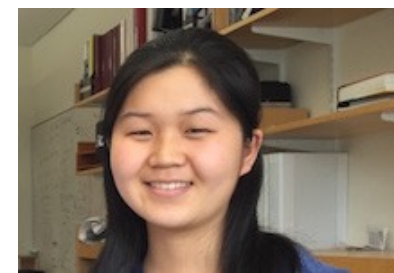What task(s)?

a) vanilla categorization

b) time-discounting

$$L = \sum_t \gamma^t \cdot L_t$$

*be accurate but also fast*

c) heavy occlusion &c

**bypass** *concat*

**memory**
$state_{t+1} = conv_{t+1} + \lambda \cdot conv_t$

*impulse transfer function*

duration: 1, 2, 3, 4, 5

time

fc, fc, fc, conv, conv, conv, conv, conv, retina, image

**Aran Nayebi** — Surya Ganguli — **Jonas Kubilius** — Maryann Rui — **Kohitij Kar** — Jim DiCarlo

fc

fc

fc

conv

conv

**bypass**
*concat*

conv

conv

conv

**memory**
$state_{t+1} = conv_{t+1} + \lambda \cdot conv_t$

**impulse transfer**
**function**

retina

image

duration
1
2
3
4
5

**time**

What task(s)?

a) vanilla categorization

b) time-discounting

$$L = \sum_t \gamma^t \cdot L_t$$

*be accurate but also fast*

c) heavy occlusion &c

strategy: optimize for tasks check against static & dynamic data

**Aran**
**Nayebi**

Surya
Ganguli

**Jonas**
**Kubilius**

Maryann
Rui

**Kohitij**
**Kar**

Jim
DiCarlo

*Petersen, 2007*

Chengxu Zhuang

Mitra Hartmann & Lab

*Petersen, 2007*



✳ Spatially-structured input data

Chengxu Zhuang

Mitra Hartmann & Lab

# Rodent Somatosensory Cortex

*Petersen, 2007*



A. From Whisker to Cortex
B. Whiskers and Barrels
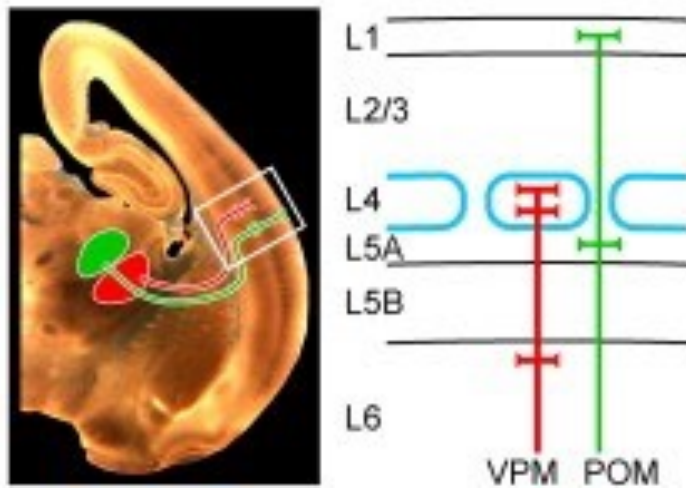C. Thalamocortical connectivity
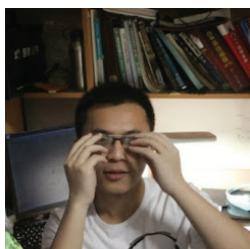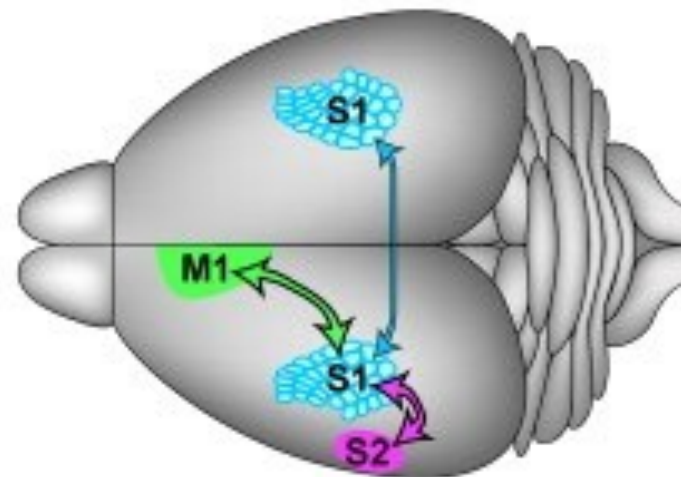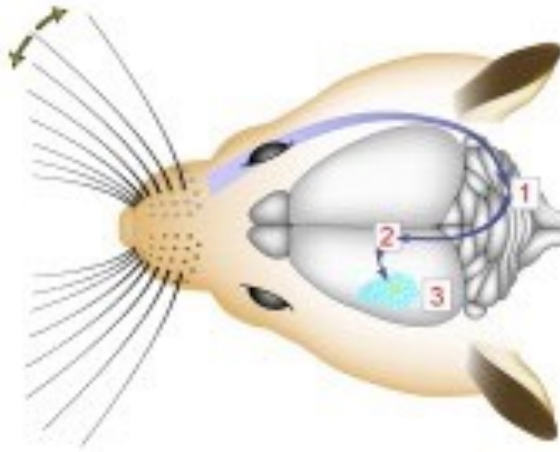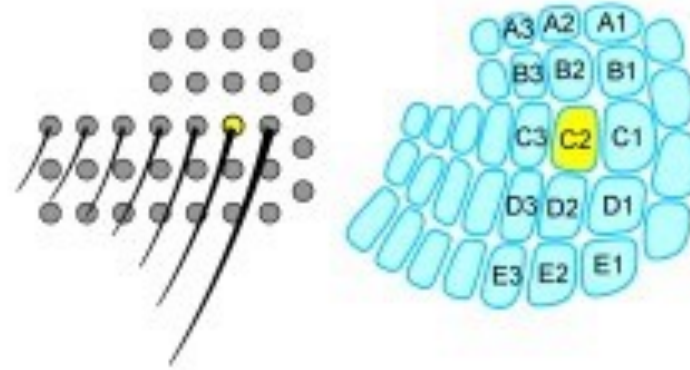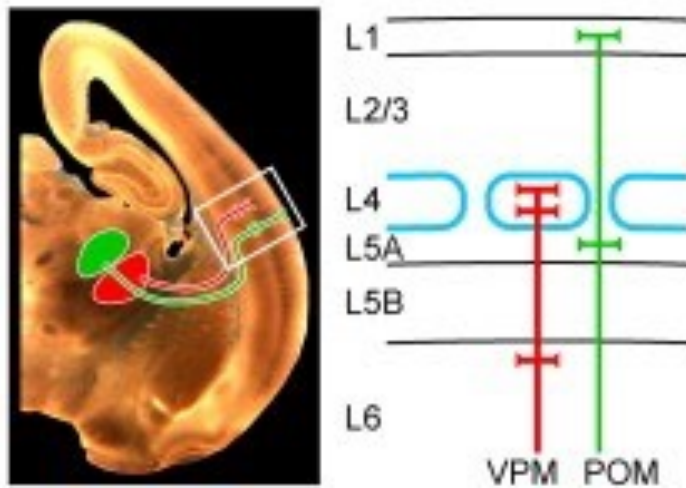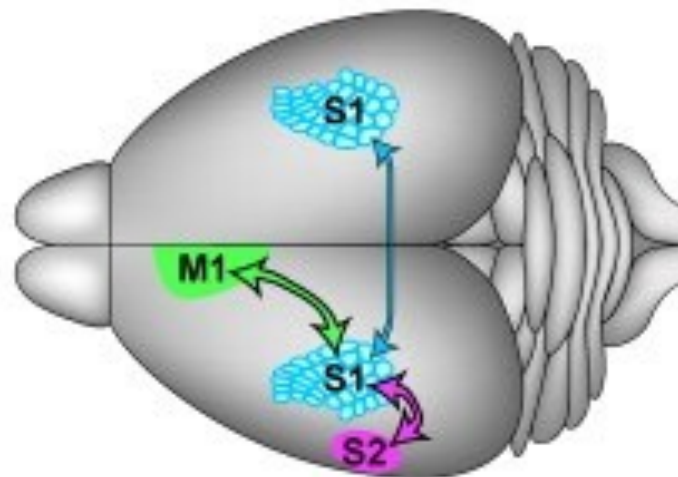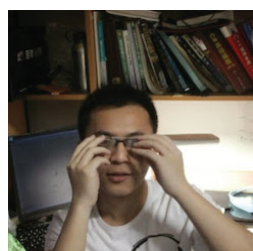D. Corticocortical connectivity

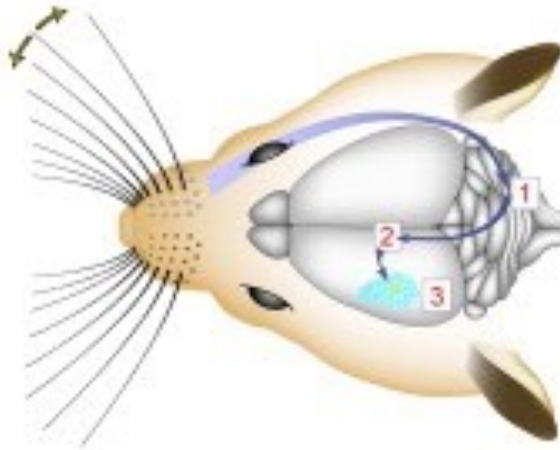* Spatially-structured input data

* Spatiotopic sensor

Chengxu Zhuang     Mitra Hartmann & Lab
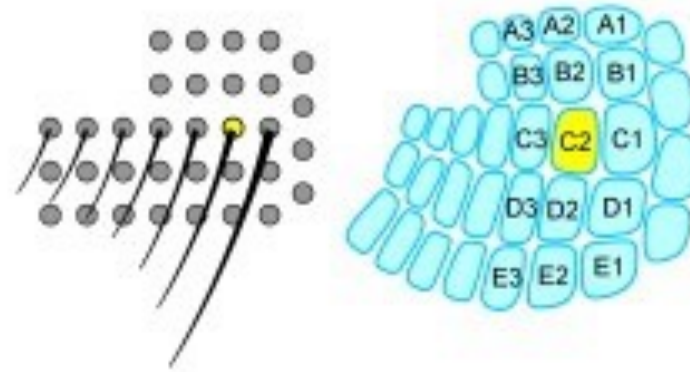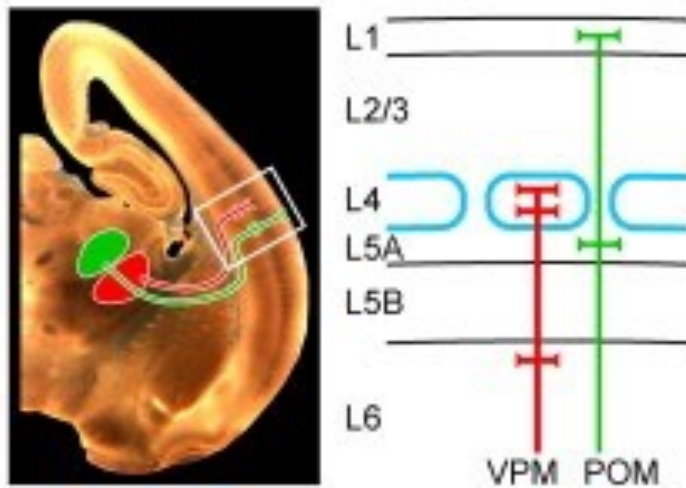
# Rodent Somatosensory Cortex

*Petersen, 2007*



A  From Whisker to Cortex
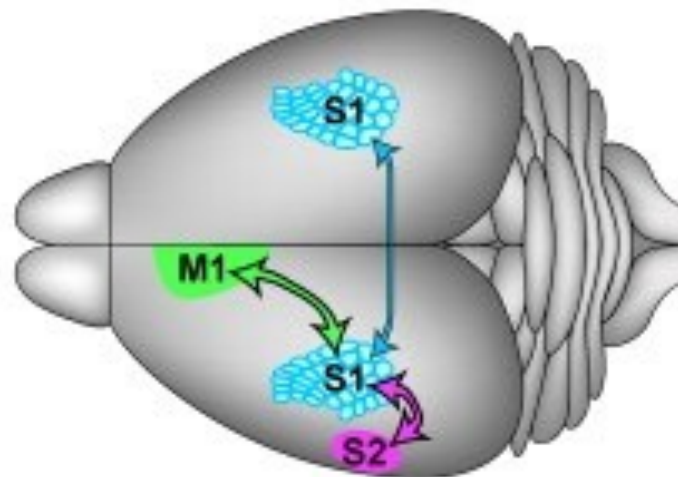
B  Whiskers and Barrels
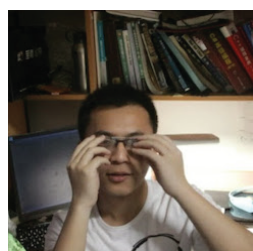
C  Thalamocortical connectivity

D  Corticocortical connectivity

* Spatially-structured input data

* Spatiotopic sensor

* Potentially hierarchical structure

Chengxu Zhuang

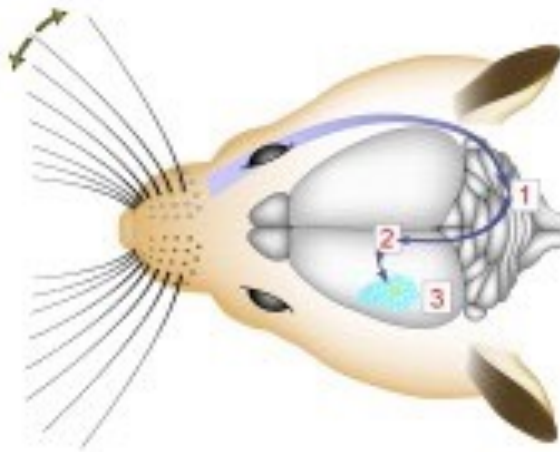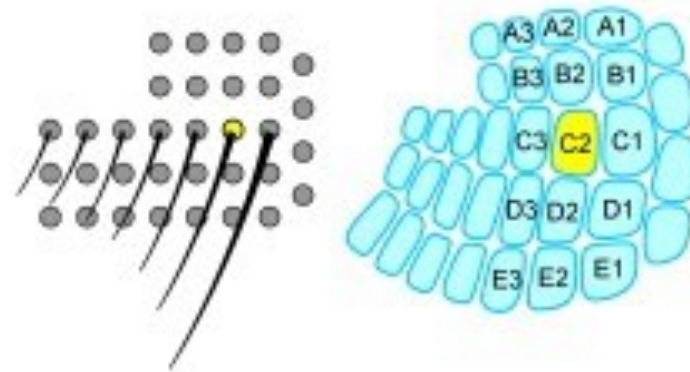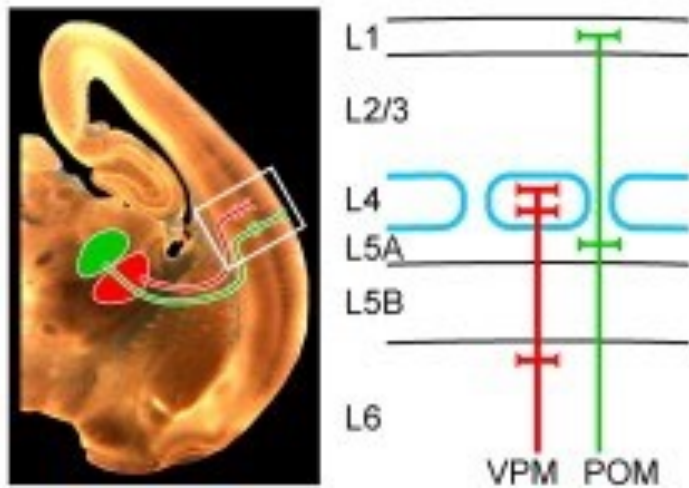Mitra Hartmann & Lab

# Rodent Somatosensory Cortex

*Petersen, 2007*



* Spatially-structured input data

* Spatiotopic sensor

* Potentially hierarchical structure

* Poorly understood higher cortical areas

Chengxu Zhuang

Mitra Hartmann & Lab

# Rodent Somatosensory Cortex



*Petersen, 2007*

* Spatially-structured input data

* Spatiotopic sensor

* Potentially hierarchical structure

* Poorly understood higher cortical areas

Chengxu Zhuang

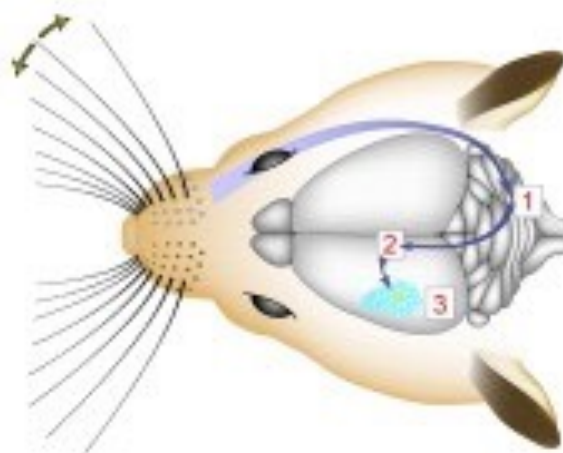Mitra Hartmann & Lab

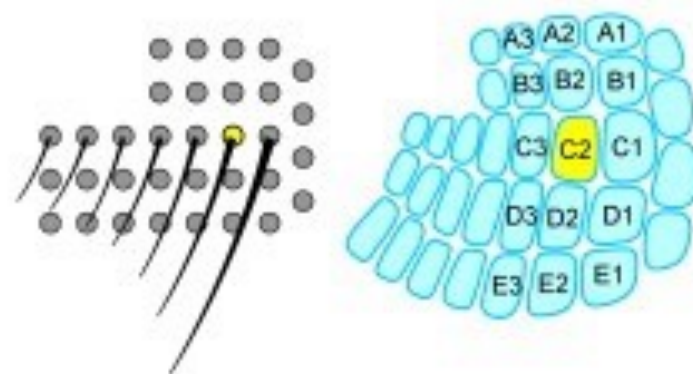Hypothesis: can get a model for this cortical cascade by optimizing properly-sized CNN with whisker-like sensor input for some ethologically relevant somatosensory task.

First have to build a model of the sensory to gather data.



Using published data from Mitra Hartmann's group



Chengxu
Zhuang

Mitra Hartmann
& Lab

First have to build a model of the sensory to gather data.



Follicle

Using published data from Mitra Hartmann's group

Chengxu
Zhuang

Mitra Hartmann
& Lab

First have to build a model of the sensory to gather data.



Chengxu
Zhuang

Mitra Hartmann
& Lab

Validate sensor on one-object tasks … (teddy vs. duck)



Testing perf vs. training num



Testing perf vs. Variation



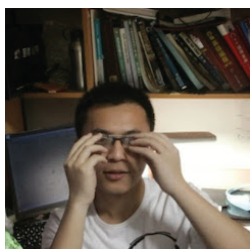Testing perf vs. Feature num

train/test splits with different:

    ✳ attack vectors

    ✳ attack speed

    ✳ object rotations

    ✳ object size



Chengxu
Zhuang

Mitra Hartmann
& Lab

Validate sensor on one-object tasks … (teddy vs. duck)



train/test splits with different:

* attack vectors

* attack speed

* object rotations

* object size

Train on shape recognition and/or normal estimation task, compare to neural data

Chengxu Zhuang

Mitra Hartmann & Lab

If successful:



Model Architecture Class

Localization

Deeper networks

Categorization

Word recognition

V1  V4  V2  PIT  CIT  AIT

??

??

??

??

Primary
Belt
Parabelt

somatosensation

vision

audition

> Formulate comprehensive model class (**CNNs**)

> Choose challenging, ethologically-valid tasks (**categorization**)

> Implement generic learning rules (**gradient descent**)

> Map to brain data. (**ventral stream**)

> Formulate comprehensive model class (**CNNs + RNs**)

> Choose challenging, ethologically-valid tasks (**task switching/ memory**)

> Implement generic learning **and expansion** rules

> Map to brain data. (**PFC, Hippocampus, &c**)



Model Architecture Class

Localization

Categorization

Deeper networks

Word recognition

Kevin Feigelis

Mark Schnitzer

Jim DiCarlo

# Q&A

Task-driven modeling can make greatly improved quantitative models of high-level cortical areas.

These models can lead to new qualitative insight about how the brain solves sensory tasks.

These concepts are useful across multiple sensory modalities.

# IT Neurons Track Human Performance

IT matches human error patterns as well as raw performance.



IT Population

V4 Population

Human Performance

Neural Decode Performance

Human Dprime

Neural D-prime (128 units)

● Low-Variation Face subordinate tasks.

# Some kind of mapping is necessary.

# Neural Response Prediction

Here, we use linear regression.

**Neural predictivity**: the ability of model to predict each individual neural site's activity.

Neural site unit ~ sparse linear combination of model units

Linear regression with fixed training images.

Accuracy = goodness-of-fit on held-out testing images (Cross validated)

Neural predictivity = median accuracy over all units.



Neural Recordings from IT and V4

**Neural predictivity**: the ability of model to predict each individual neural site's activity.

Neural Recordings from IT and V4

# Performance Comparison



Yamins* and Hong* et. al. **PNAS** (2014)

Low Var.     Medium Var.     High Var.

Basic categorization

Car identification

Face identification

Performance

Pixels, SIFT, V1-like, V2-like, HMAX, PLOS09, V4 Neurons, IT Neurons, Human, HMO

# Neural Data Recording

*Output = Binned spike counts in 70ms-170ms post stimulus presentation; averaged over 25-50 reps of each image.*

# Single Site Responses



Site 10    Site 54    Site 43

Site 11    Site 77    Site 102

Best single position-encoding sites.

heat map value at x, y =
 response averaged over all
 images where object center is in
 position x, y

# Single Site Responses

Site 10          Site 54          Site 43

*y*

*x*

Site 11          Site 77          Site 102

Best single position-encoding sites.

heat map value at x, y =
    response averaged over all
    images where object center is in
    position x, y

Similar to MacEvoy (2013) and DiCarlo(2003)
except — dramatically more variation.

# Single Site Responses



Best single position-encoding sites.

heat map value at x, y =
response averaged over all
images where object center is in
position x, y

# Monkey Neurons vs Humans

# Monkey Neurons vs Humans

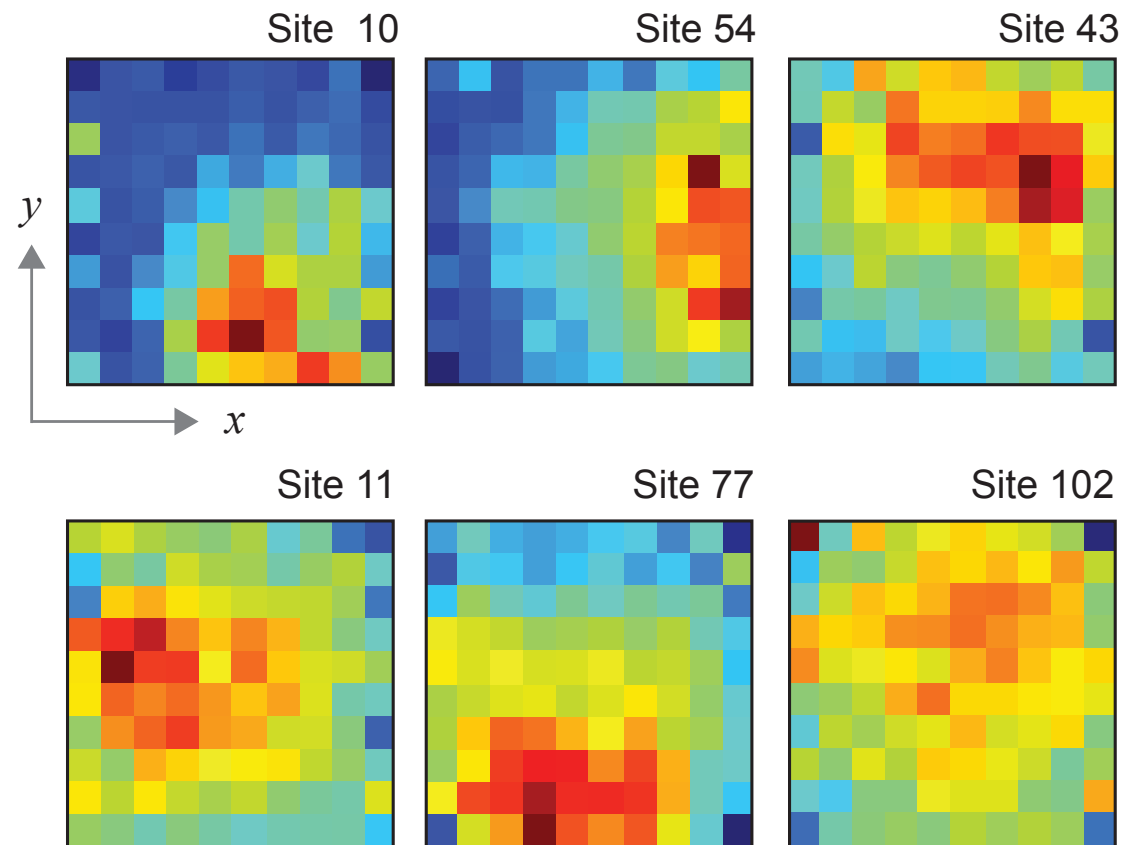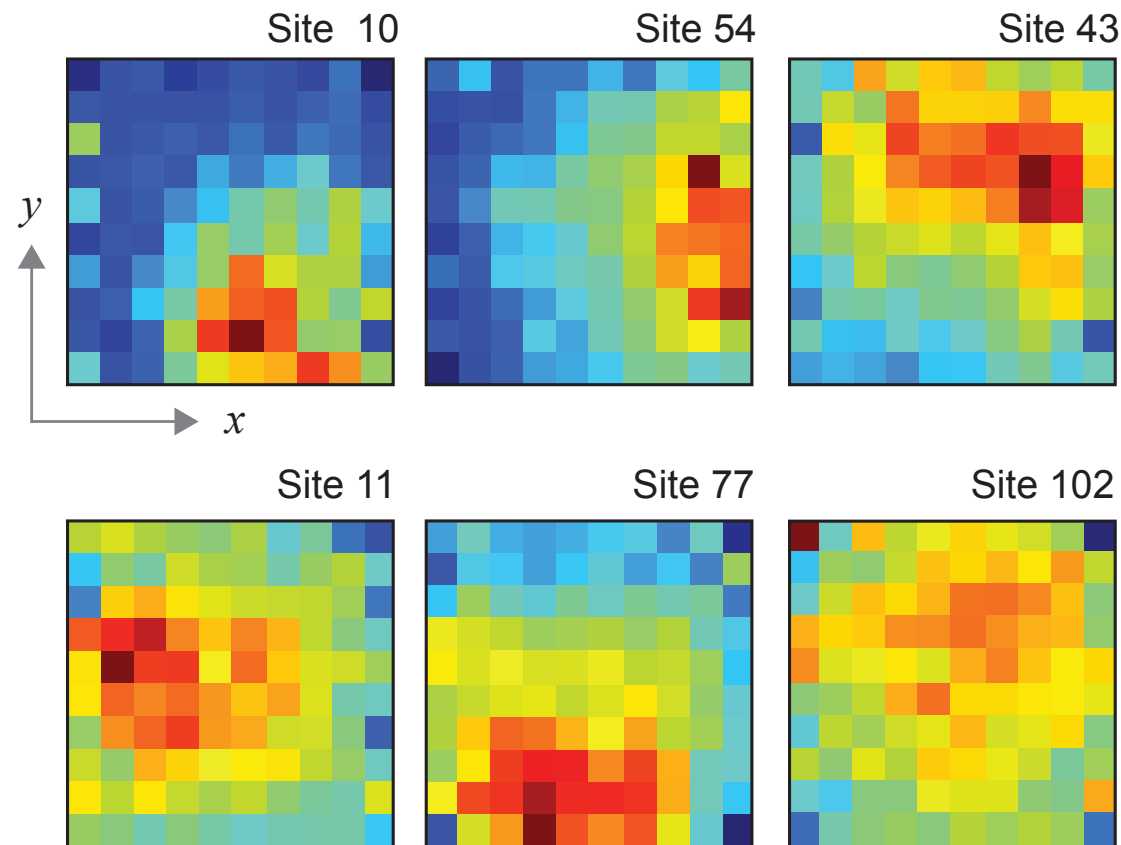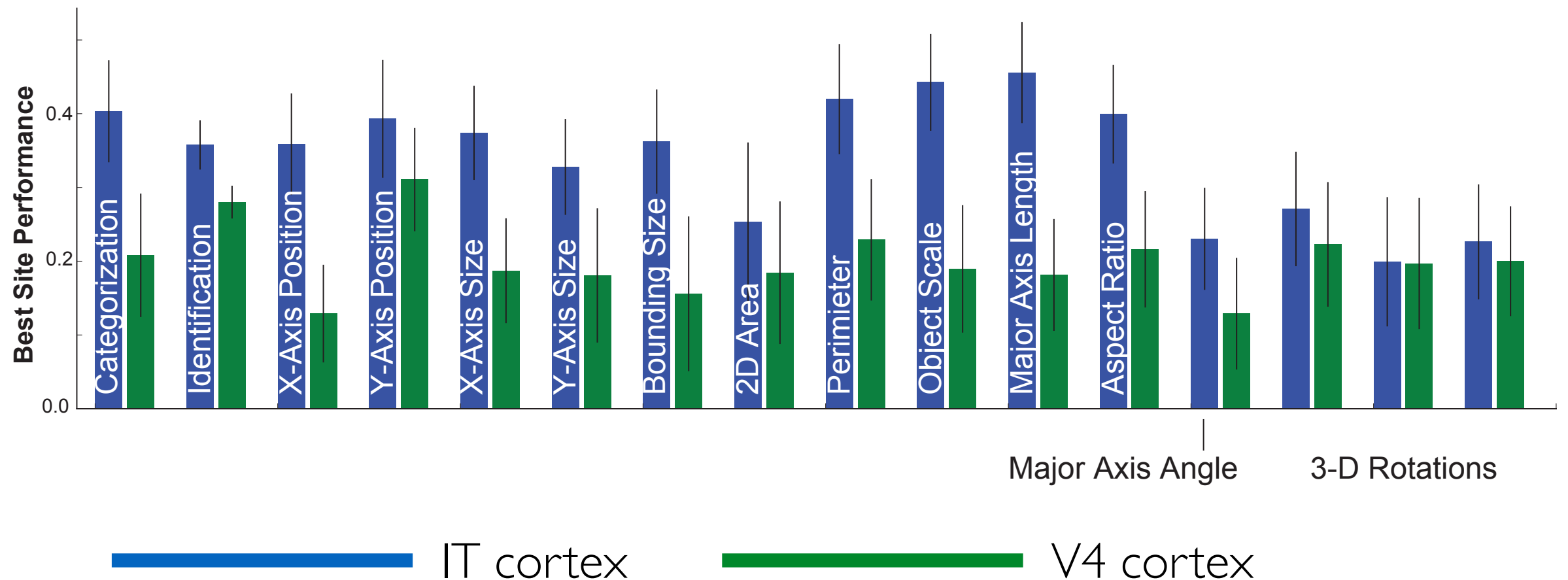| | IT | V4 | V1 | Pix |
|---|---|---|---|---|
| Basic Categorization | 773 ± 185 | $2.2 \times 10^6$ | — | — |
| Subordinate Identification | 496 ± 93 | $4.4 \times 10^6$ | — | — |
| X-axis Position | 1414 ± 403 | $5.2 \times 10^5$ | $3.0 \times 10^7$ | — |
| Y-axis Position | 918 ± 309 | $2.5 \times 10^4$ | $8.7 \times 10^6$ | — |
| Bounding Box Size | 322 ± 90 | $1.7 \times 10^4$ | — | — |
| X-axis Size | 256 ± 87 | $9.8 \times 10^3$ | $3.4 \times 10^7$ | — |
| Y-axis Size | 237 ± 87 | $3.8 \times 10^3$ | $9.5 \times 10^6$ | — |
| 3-D Object Scale | 401 ± 90 | $3.2 \times 10^4$ | — | — |
| Major Axis Length | 201 ± 70 | $1.1 \times 10^4$ | — | — |
| Aspect Ratio | 163 ± 61 | 951 ± 59 | $6.5 \times 10^3$ | — |
| Major Axis Angle | 804 ± 136 | $3.2 \times 10^6$ | — | — |
| Z-axis Rotation | 1932 ± 1061 | — | — | — |
| Y-axis Rotation | 369 ± 115 | $2.8 \times 10^5$ | — | — |
| X-axis Rotation | 1570 ± 530 | — | — | — |

— = more than 10 billion sites required

*Hong\*, Yamins\*, Majaj & DiCarlo.* **Nat. Neuro.** *(in press)*

Mean over tasks, human-parity for IT is at ~**700** multi-unit trial-averaged sites.

# Monkey Neurons vs Humans

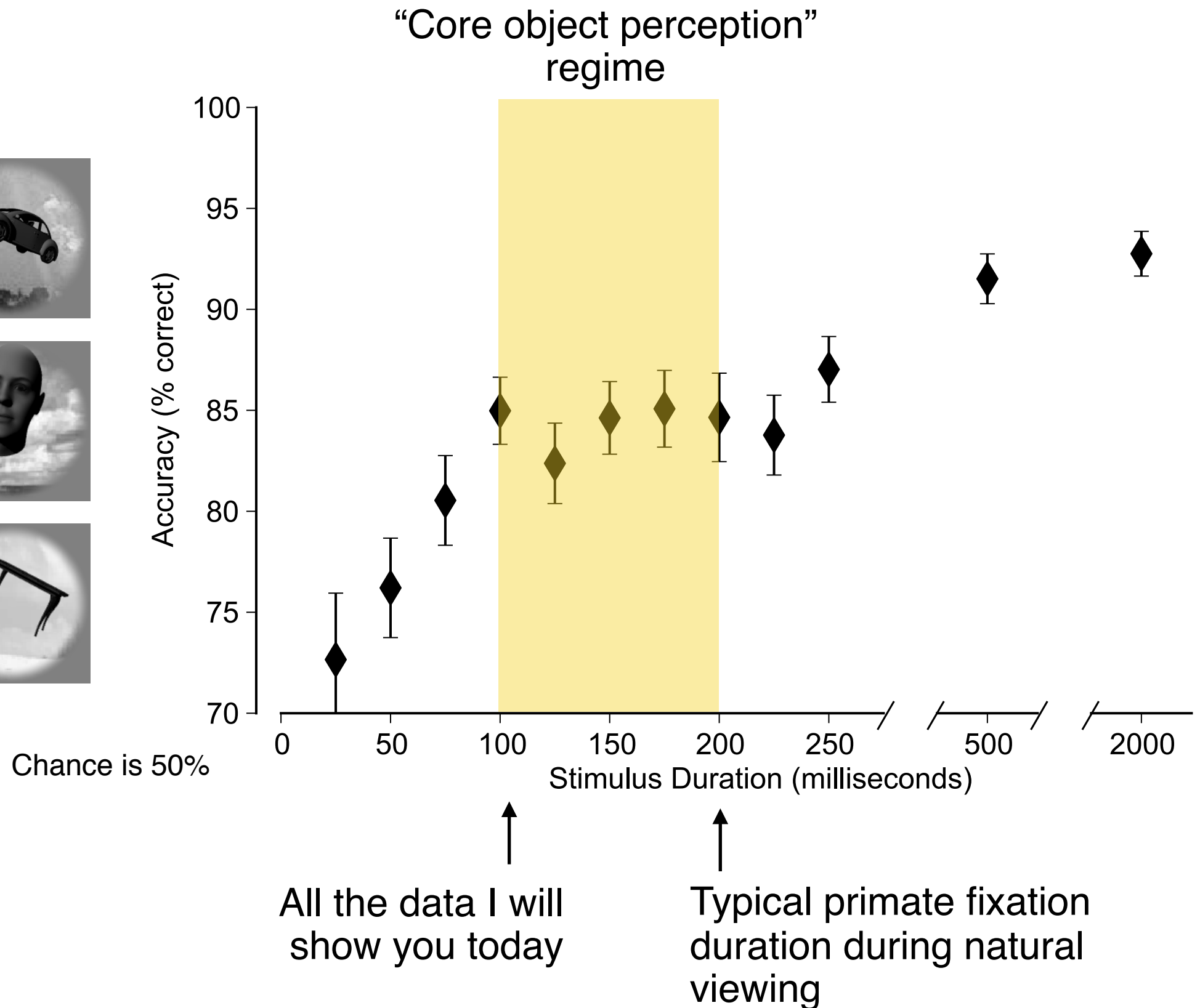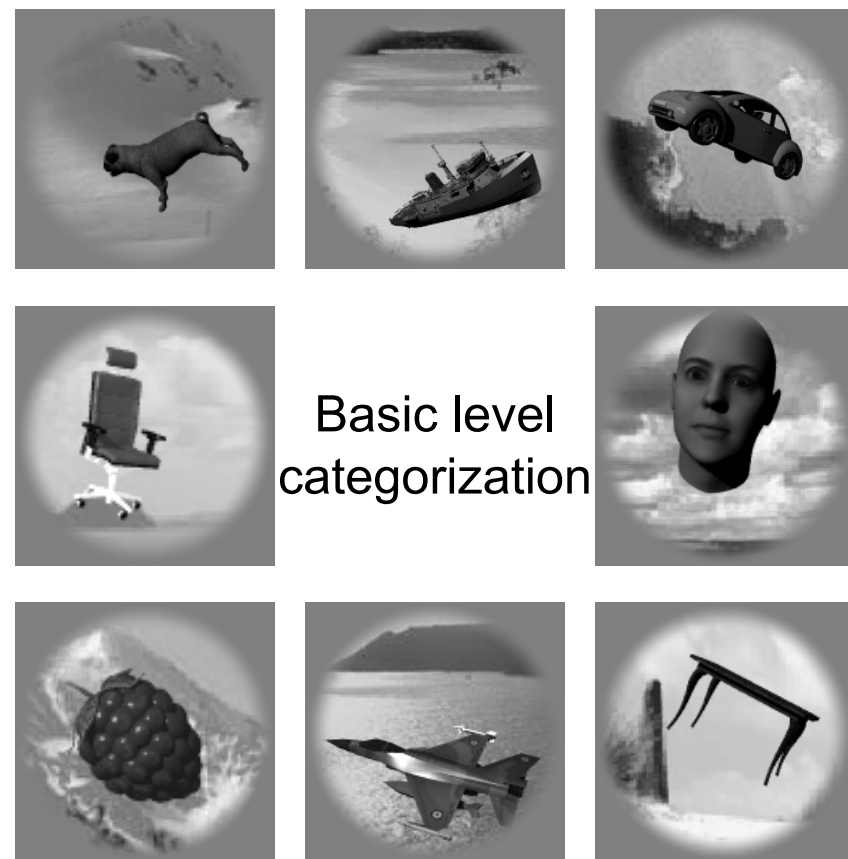|  | IT | V4 | V1 | Pix |
|---|---|---|---|---|
| Basic Categorization | 773 ± 185 | $2.2 \times 10^6$ | — | — |
| Subordinate Identification | 496 ± 93 | $4.4 \times 10^6$ | — | — |
| X-axis Position | 1414 ± 403 | $5.2 \times 10^5$ | $3.0 \times 10^7$ | — |
| Y-axis Position | 918 ± 309 | $2.5 \times 10^4$ | $8.7 \times 10^6$ | — |
| Bounding Box Size | 322 ± 90 | $1.7 \times 10^4$ | — | — |
| X-axis Size | 256 ± 87 | $9.8 \times 10^3$ | $3.4 \times 10^7$ | — |
| Y-axis Size | 237 ± 87 | $3.8 \times 10^3$ | $9.5 \times 10^6$ | — |
| 3-D Object Scale | 401 ± 90 | $3.2 \times 10^4$ | — | — |
| Major Axis Length | 201 ± 70 | $1.1 \times 10^4$ | — | — |
| Aspect Ratio | 163 ± 61 | 951 ± 59 | $6.5 \times 10^3$ | — |
| Major Axis Angle | 804 ± 136 | $3.2 \times 10^6$ | — | — |
| Z-axis Rotation | 1932 ± 1061 | — | — | — |
| Y-axis Rotation | 369 ± 115 | $2.8 \times 10^5$ | — | — |
| X-axis Rotation | 1570 ± 530 | — | — | — |

— = more than 10 billion sites required

*Hong\*, Yamins\*, Majaj & DiCarlo.* **Nat. Neuro.** *(in press)*

Mean over tasks, human-parity for IT is at ~**350000** single-unit single-trial neurons.

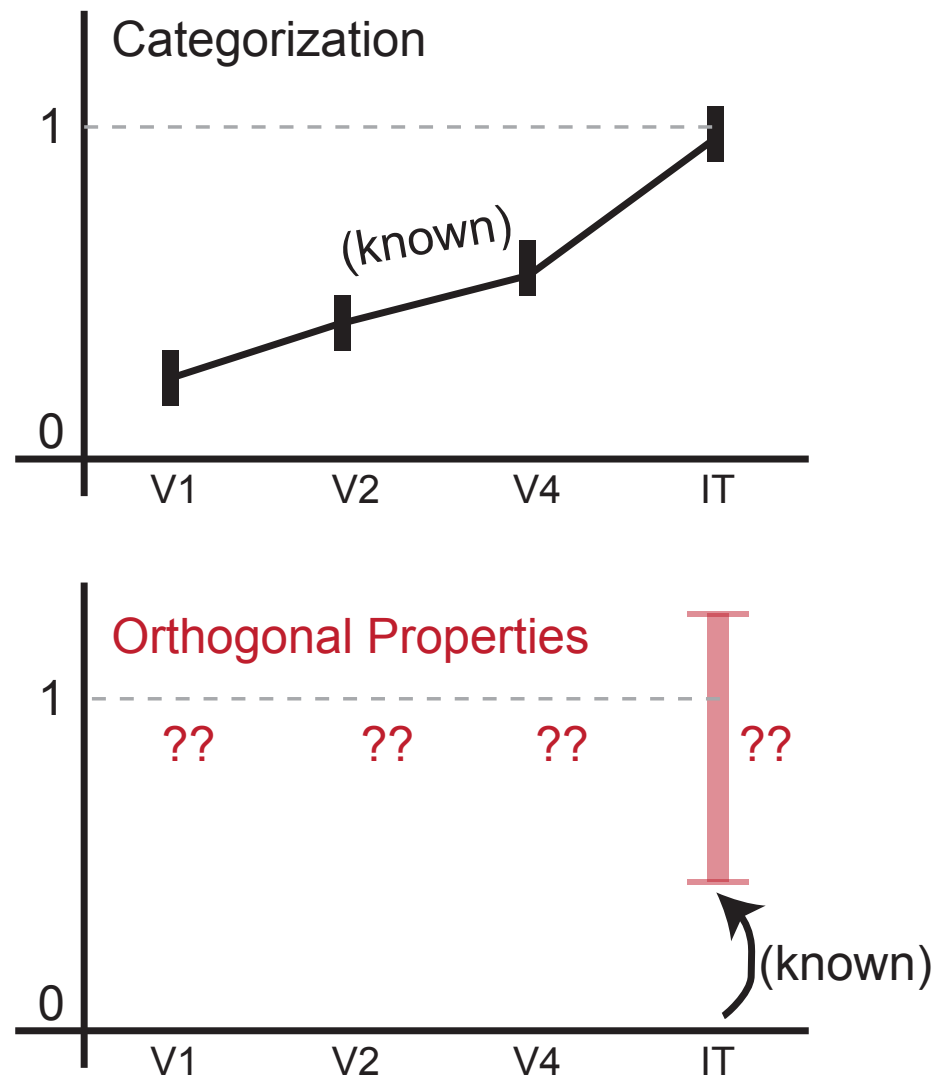# Example: Human object categorization accuracy as a function of image viewing time



Basic level categorization

Chance is 50%

"Core object perception" regime

Accuracy (% correct)

Stimulus Duration (milliseconds)

All the data I will show you today

Typical primate fixation duration during natural viewing

Somewhat newish ideas about IT?

Population Decode Performance
*(relative to human performance)*
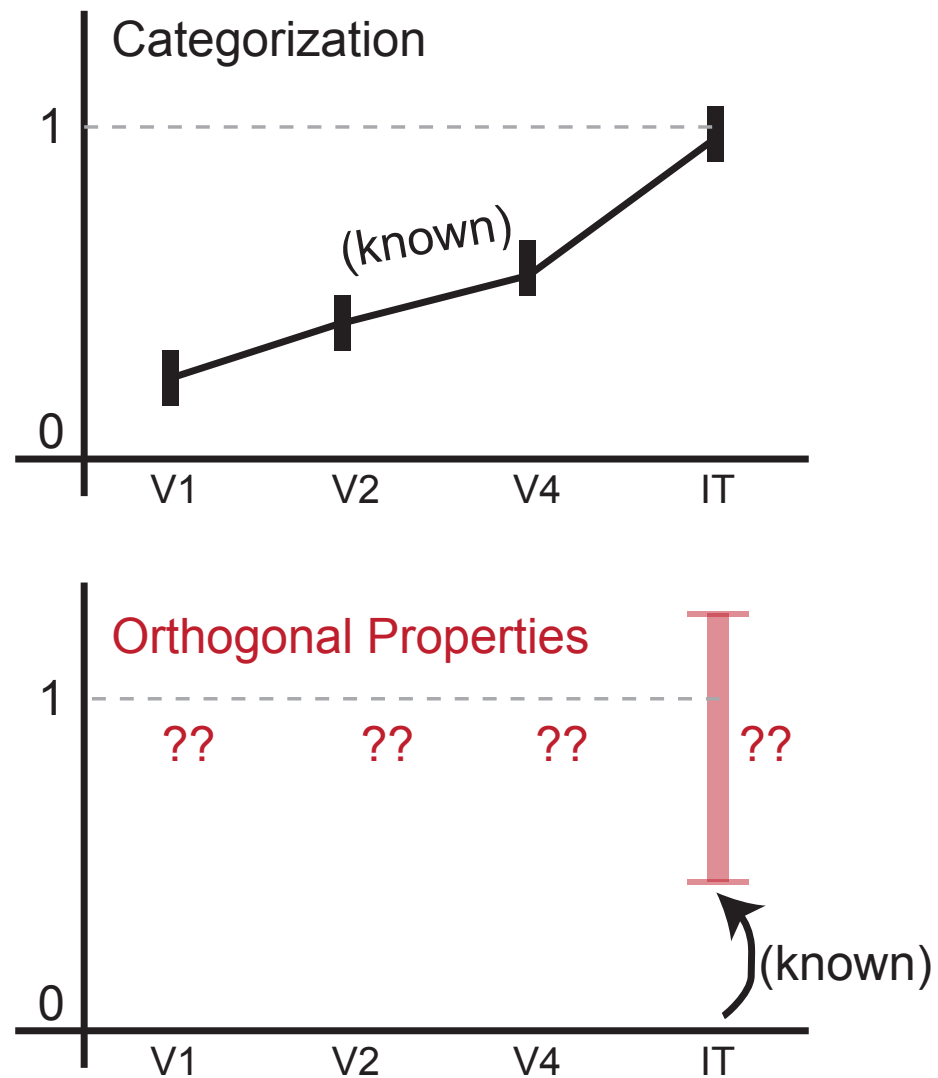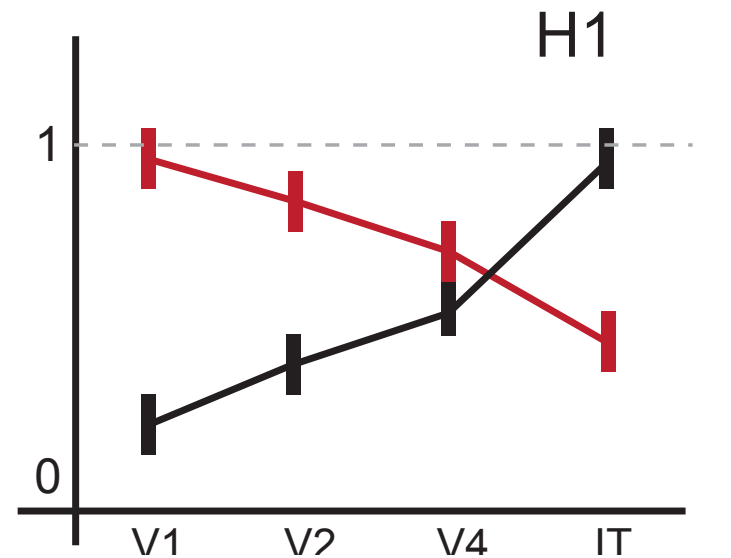
State of knowledge from previous studies . . .

Categorization

1

(known)

0

V1    V2    V4    IT

Orthogonal Properties

1    ??    ??    ??    ??

0

V1    V2    V4    IT

(known)

Multiple hypotheses consistent with the existing data . . .

H1

1

0

V1    V2    V4    IT
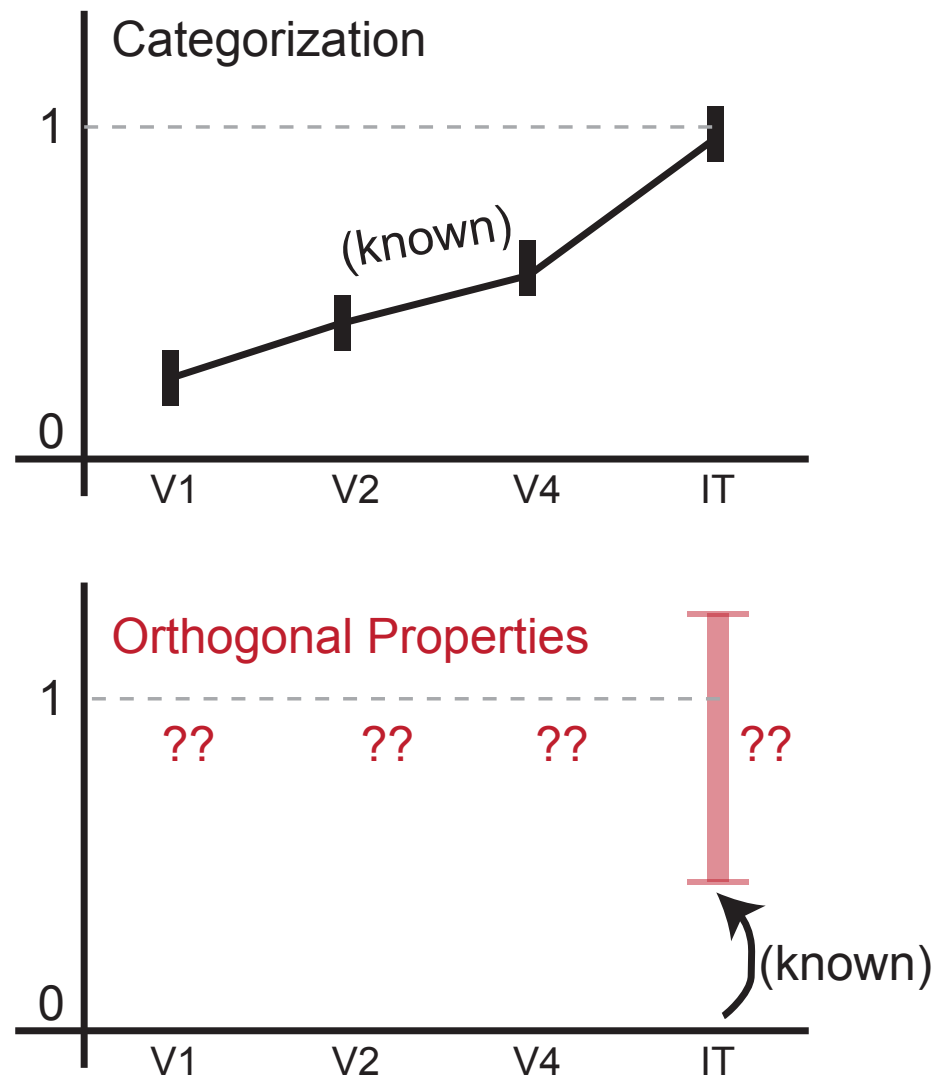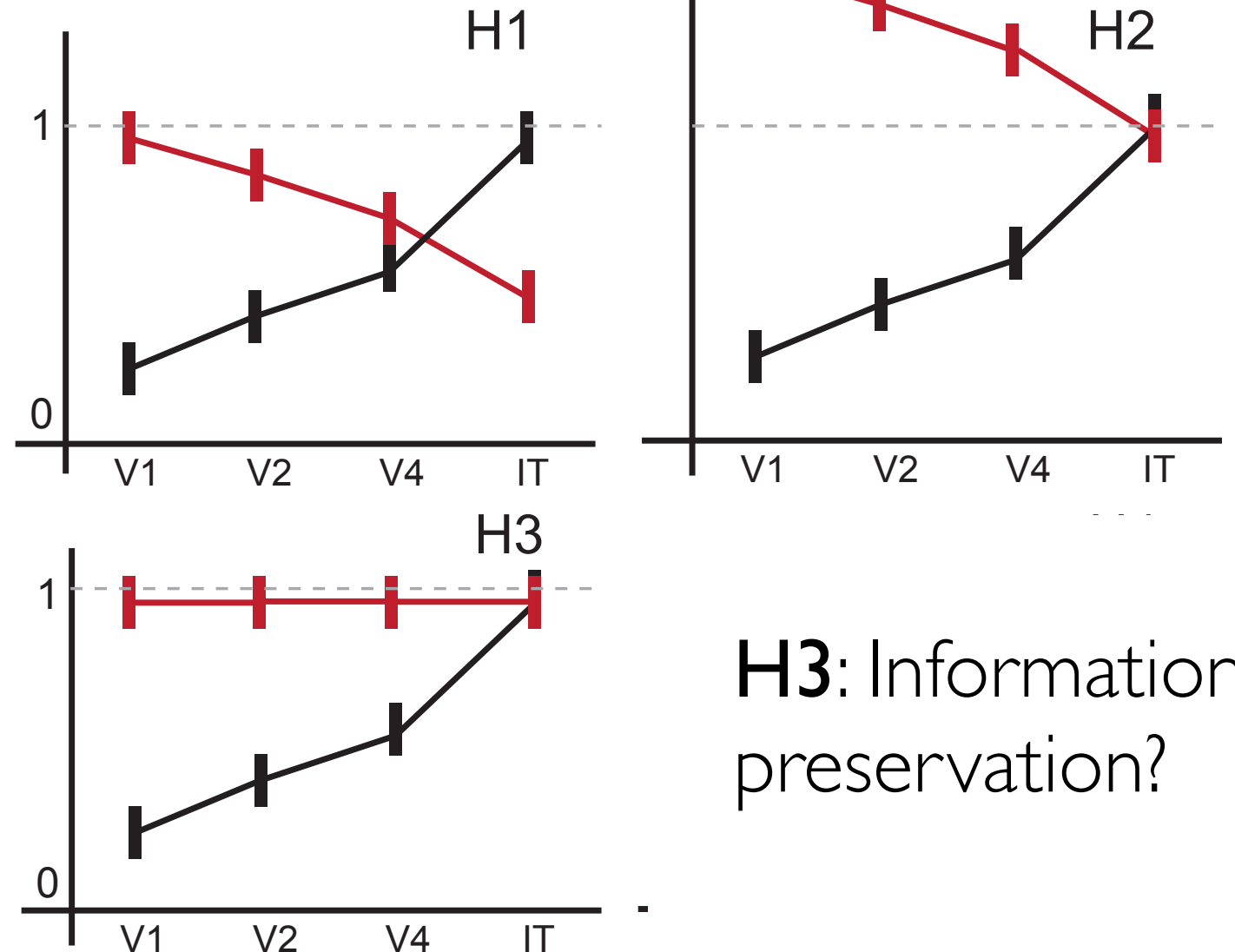
H1: Tolerance / sensitivity tradeoff?

Depth Along Ventral Stream
*(increasing receptive field size →)*

Somewhat newish ideas about IT?

Somewhat newish ideas about IT?

State of knowledge from previous studies . . .

Multiple hypotheses consistent with the existing data . . .

Population Decode Performance *(relative to human performance)*

Depth Along Ventral Stream *(increasing receptive field size →)*

**H4**: Simultaneous build-up of encoding

TEOd

TEpd

Face patches

TEad

TGv
granular

*R. Lafer-Sousa and BR Conway, Nat. Neurosci (2013)*

Regions selective for:

- faces

TEOd

TEpd

TEad

TGv
granular

Face patches

Color-biased regions

*R. Lafer-Sousa and BR Conway, Nat. Neurosci (2013)*

Regions selective for:
- faces
- places
- bodies
- color

# Selective Patches in Higher Visual Cortex



TEOd
TEpd
TEad
TGv granular

*R. Lafer-Sousa and BR Conway, Nat. Neurosci (2013)*

Face patches

Color-biased regions

Regions selective for:

- faces

- places

- bodies

- color

Where do these patches come from?

- In-born built-in structure??

- or developmentally determined by domain-specific experience?

controlled rearing

controlled rearing



−8
−3

−8
−2

????

TEOd

TEpd

TEad

TGv
granular

Face patches
Color-biased regions

controlled rearing

TEOd

TEpd

TEad

TGv
granular

−8
−3

−8
−2

Face patches
Color-biased regions

????

… in a computational model

ImageNet (2012). Thousands of images in thousands of categories.

controlled rearing



remove all images containing faces
as well as all categories of photos containing animate objects

**question:** how does removing this content affect the model?

revolver

cheeseburger

tiger

german shepherd

bolete

sycamore

laptop

beeker

television

catamaran

thatched roof

parakeet

revolver

cheeseburger

tiger

german shepherd

bolete

sycamore

laptop

beeker

television

catamaran

thatched roof

parakeet

revolver

cheeseburger

tiger

german shepherd

bolete

sycamore

laptop

beeker

television

catamaran

thatched roof

parakeet

# Testing the base-line non-face model.

Boats

Cars

Chairs

⋮

Fruits

Planes

Tables

① *Hierarchical Convolutional Neural Network ...*

**layer 1**

② *... trained on images with no faces or animate objects ...*

**layer 6**

③ *... face-selective units identified with a standard localizer ...*

Original speculation: we won't find any (or statistically significantly many) face selective units because:

Hypotheses for existing of face-selective units:

  i. ~~face processing machinery is in-born~~

Original speculation:  we won't find any (or statistically significantly many) face selective units because:

Hypotheses for existing of face-selective units:

~~i.   face processing machinery is in-born~~

~~ii.  it is due to extensive post-natal experience with faces~~

*face selective* ⇒ d' vs non-face > 1



*Kanwisher, 1997*



Fraction Selective Units
*(d' vs all > 1)*

0.07

0.01

Faces

Testing the base-line non-face model.

*face selective* ⇒ d' vs non-face > 1

Fraction Selective Units
*(d' vs all > 1)*

***

0.07

0.01

Faces

…

~7% of units in model were face selective

*face selective* ⇒  d' vs non-face > 1



Fraction Selective Units
*(d' vs all > 1)*

\*\*\*

0.07

\*

ns        ns

0.01

Faces      Cars      Planes      Tables

...  →

lower (or n.s.)
numbers for several
other tested categories

# Testing the base-line non-face model.

*face selective* ⇒ d' vs non-face > 1

# Validating the face-selective units

Boats

Cars

Chairs

Fruits

Planes

Tables

① *Hierarchical Convolutional Neural Network ...*

**layer 1**

② *... trained on images with no faces or animate objects ...*

**layer 6**

*screen for selectivity*

③ *... face-selective units identified with a standard localizer ...*

vs. vs.

vs. vs.

④ *... validated on a distinct set of testing images.*

# Validating the face-selective units

average ranked response over all face-selective units



averaged over images within category →

How could this result be true?

2-d MDS of three-d mesh distances for 128 objects in 16 categories.



How could this result be true?

# Possible explanation

2-d MDS of three-3 mesh distances for 128 objects in 16 categories



**Faces**
**Cars**
**Guns**
**Boats**
**Planes**
**Shoes**

1) Faces are more clustered in shape space than most other categories

2) but they're not totally isolated in shape space.

# Possible explanation

2-d MDS of three-3 mesh distances for 128 objects in 16 categories



**Faces**
**Cars**
**Guns**
**Boats**
**Planes**
**Shoes**

◯  =  unit as gaussian blob in shape space.

# Possible explanation

2-d MDS of three-3 mesh distances for 128 objects in 16 categories



**Faces** Cars **Guns** Boats Planes **Shoes**

◯  =  unit as gaussian blob in shape space.

detailed comparison to face neurons

IT Unit 53

Response Magnitude

Animals    Boats    Cars    Chairs    **Faces**    Fruits    Planes    Tables

Images sorted by **category**

—— Neural data

# Predictions of Face-Selective Neural Responses

IT Unit 53



Response Magnitude

Animals    Boats    Cars    Chairs    **Faces**    Fruits    Planes    Tables

Images sorted by **category**

—————  Neural data

Regularized linear regression to map model units to neural units,

predictions on held-out testing images.

# Predictions of Face-Selective Neural Responses

IT Unit 53

r^2=0.55

Response Magnitude

Animals   Boats   Cars   Chairs   **Faces**   Fruits   Planes   Tables

Images sorted by **category**

Neural data

Model prediction

# Predictions of Face-Selective Neural Responses

IT Unit 53

Response Magnitude

r^2=0.55

Images sorted by **category**

Animals  Boats  Cars  Chairs  **Faces**  Fruits  Planes  Tables
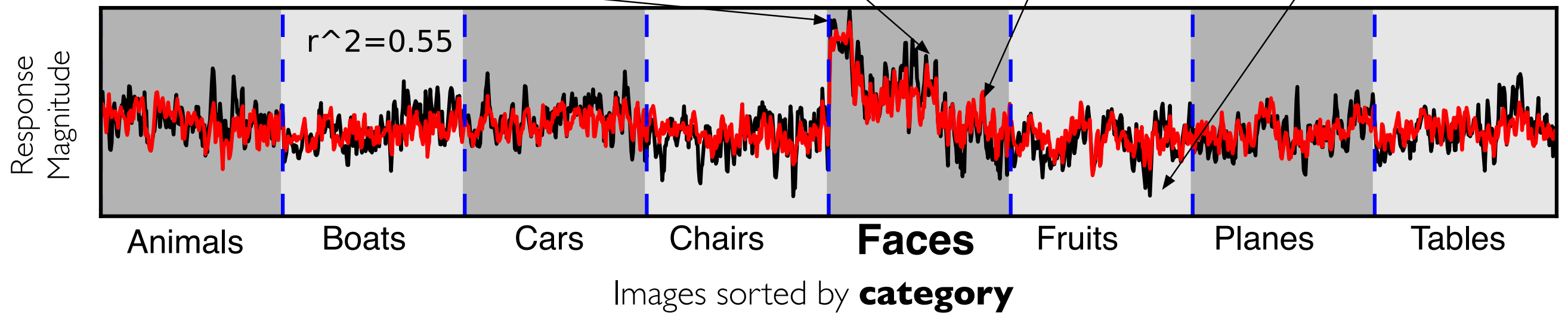
——— Neural data

——— Model prediction

Explained Variance Across All Face Selective Units:

With Faces in Training: **51.5** ± 3.9 %

Without Faces in Training: **50.8** ± 4.4 %

Models "raised" without faces can still have face-selective units.

Consistent with Sugita (2008)

Some aspects of specialized face machinery may be explicable from the "null model" of general object recognition.

A third hypothesis for the development of face (and other) selective regions:

- In-born built-in structure or

- Developmentally determined by particular experience.

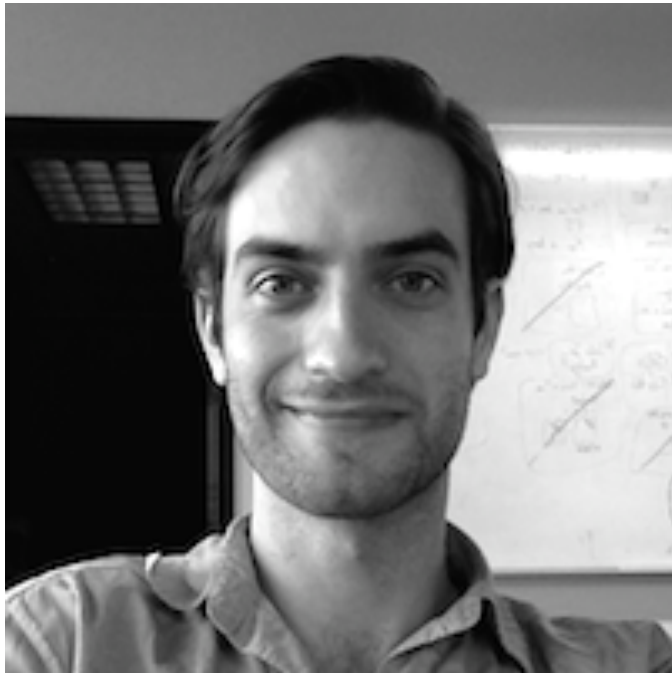- **Developmentally determined by general experience?**

Better exploration of category selective across many categories as a function of contents of training data.

More detailed comparison to neurophysiology of face patch system. Freiwald & Tsao, 2011, Issa & DiCarlo 2014
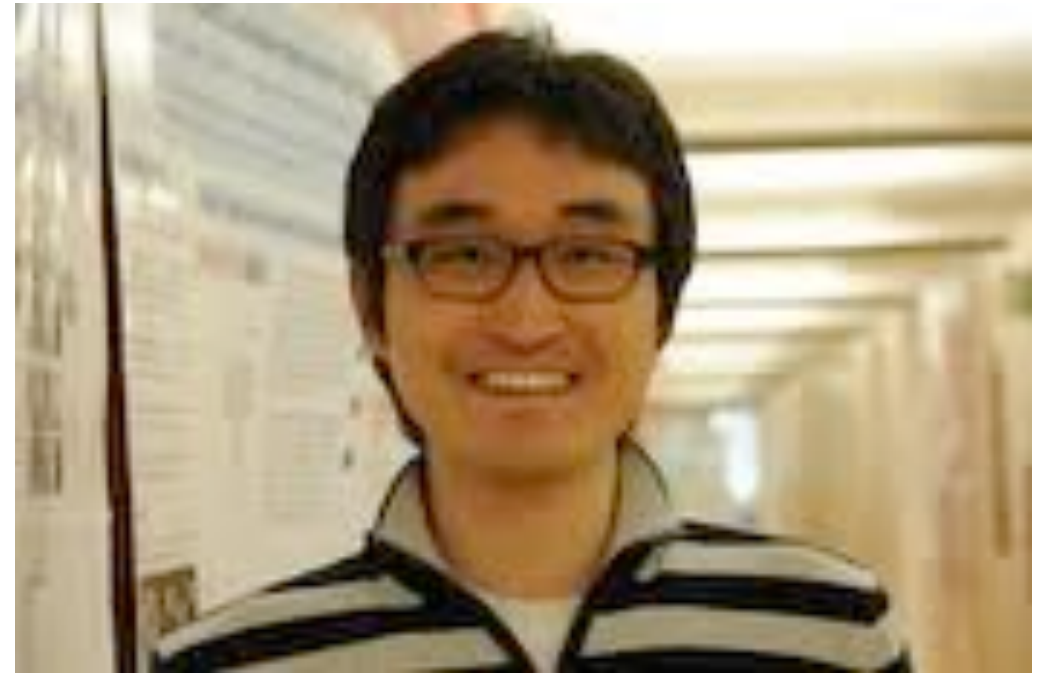
Explicitly address question of spatial layout.

Results here do <u>NOT</u> imply monkeys without face experience will necessarily have a \***patch**\*.
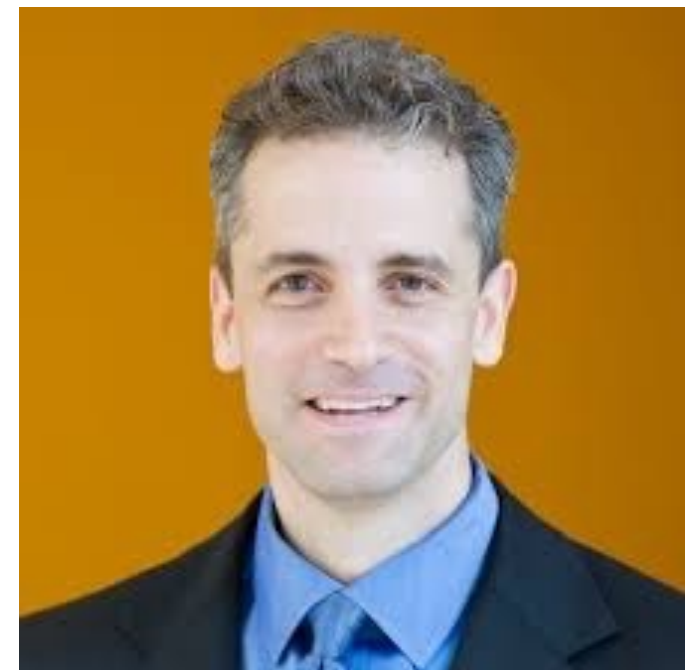
# Thanks to great colleagues!



Michael Cohen

Ha Hong

Nancy Kanwisher

Jim DiCarlo