# NX-414: Brain-like computation and intelligence

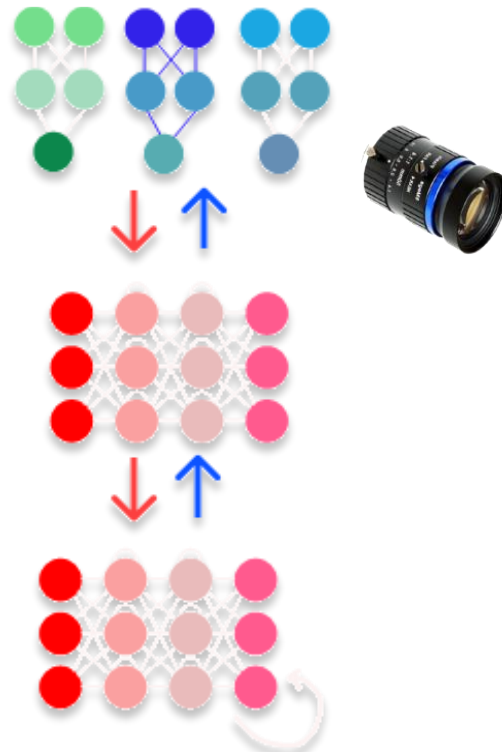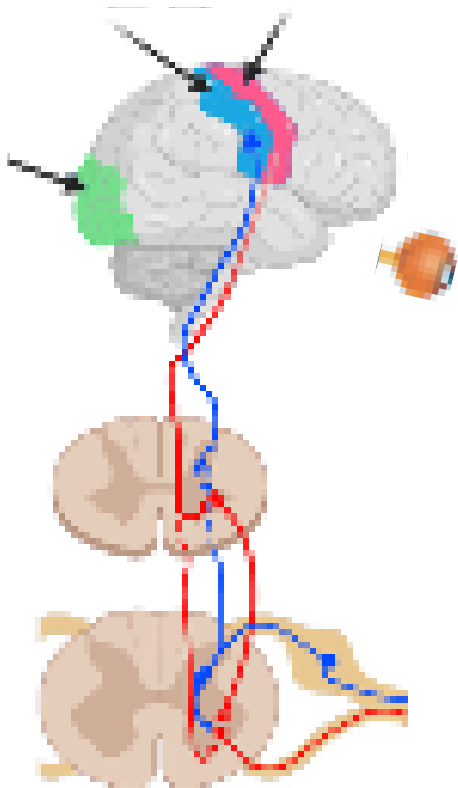Martin Schrimpf

Lecture 6, March 26 2025

# EPFL

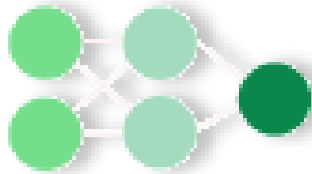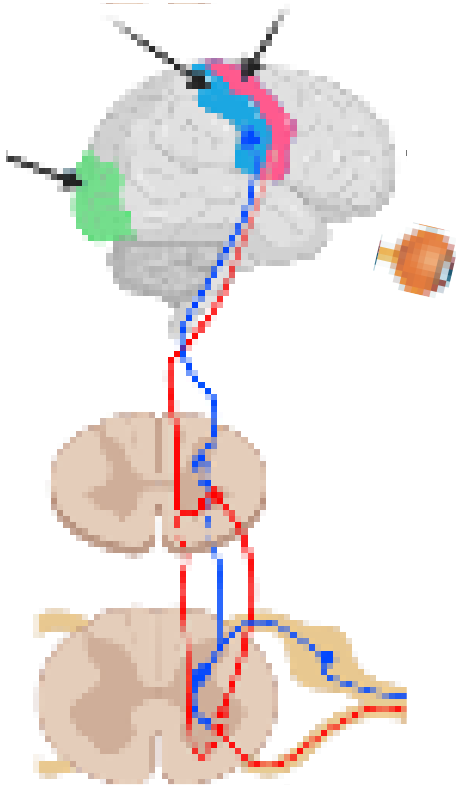## Biological Intelligence ⟷ Artificial Intelligence
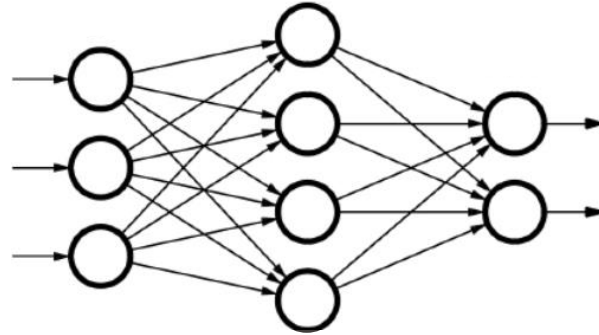
# Normative frameworks

**Information theoretic**

e.g. sparse coding, redundancy reduction, mutual information …

**Utilitarian**

e.g. **recognize objects**, chase prey, navigate …

# Using deep neural networks as goal-driven models of a system



Yamins & DiCarlo (2016)
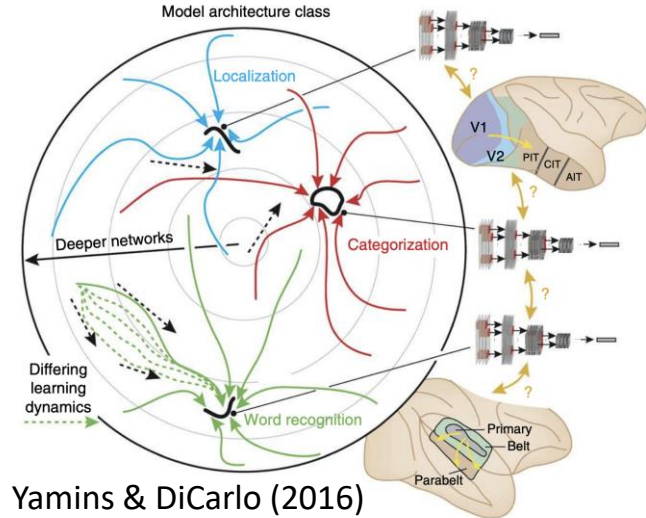
<u>Vision</u>: object recognition.
Yamins & Hong et al. (2014), Schrimpf & Kubilius et al. (2018)

<u>Audition</u>: speech recognition, speaker & sound identification. Kell et al. (2018)

<u>Somatosentation</u>: shape recognition. Zhuang et al. (2017)

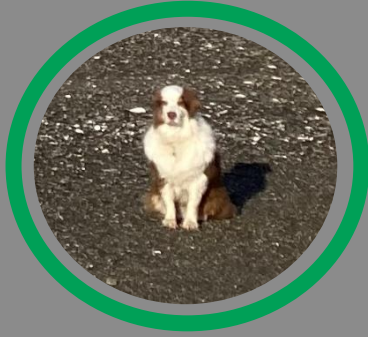<u>Language</u>: next-word prediction. Schrimpf et al. (2021)

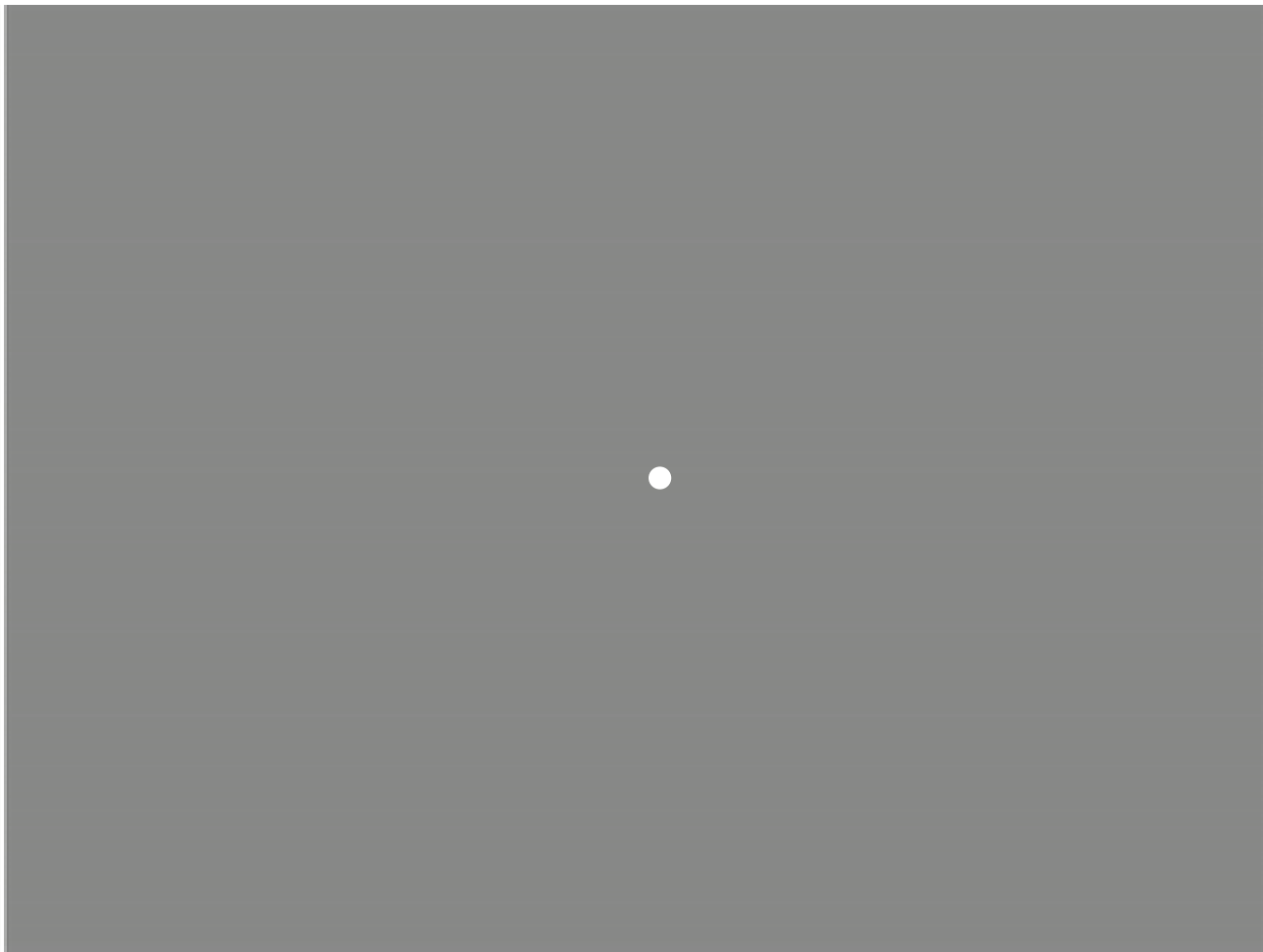<u>Decision making</u>: context-dependent choice. Mante & Sussilo et al. (2013)

<u>Proprioception</u>: action recognition. Sandbrink et al. (2023)
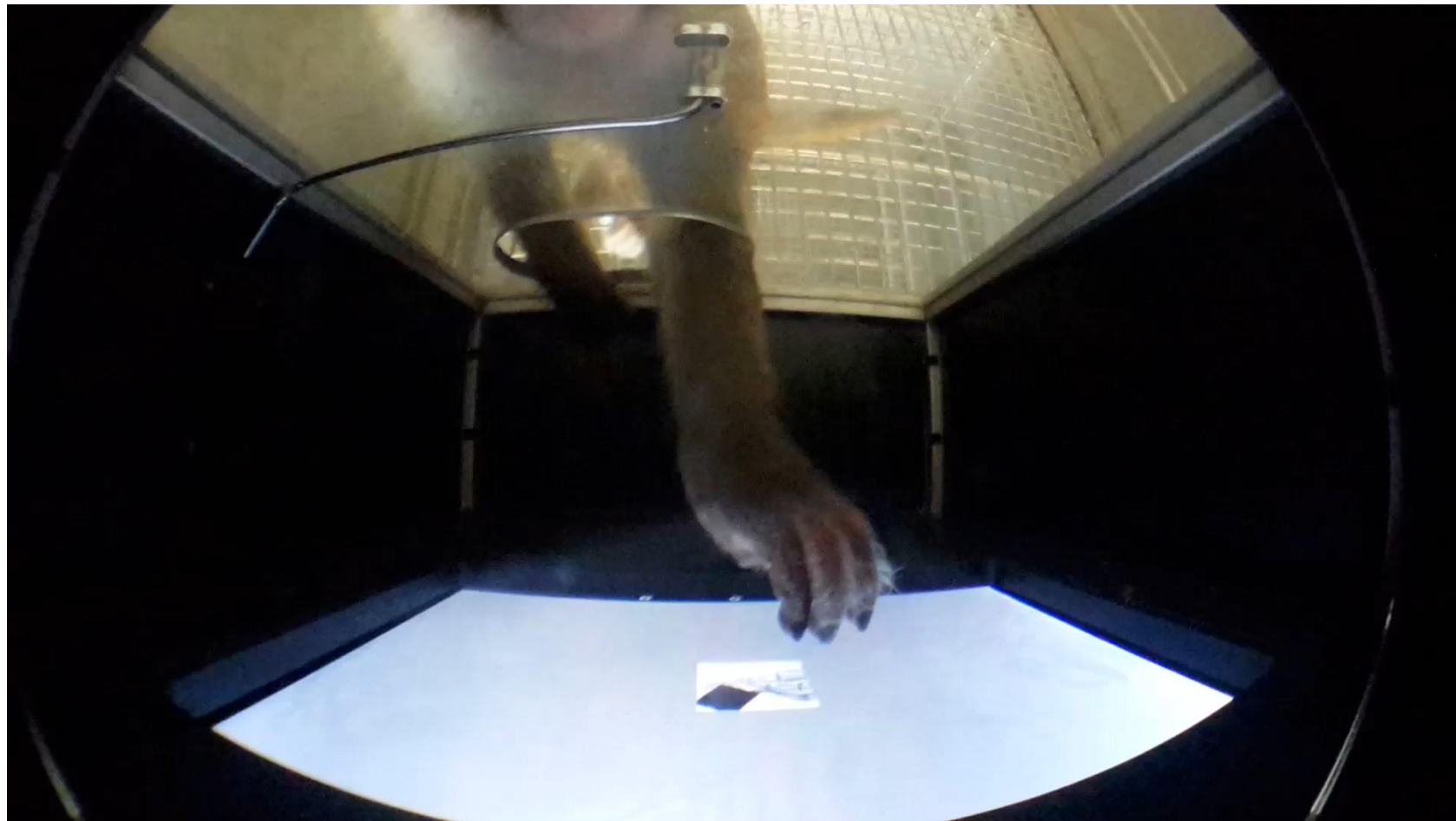
Behavioral experiment

**Behavioral experiment**
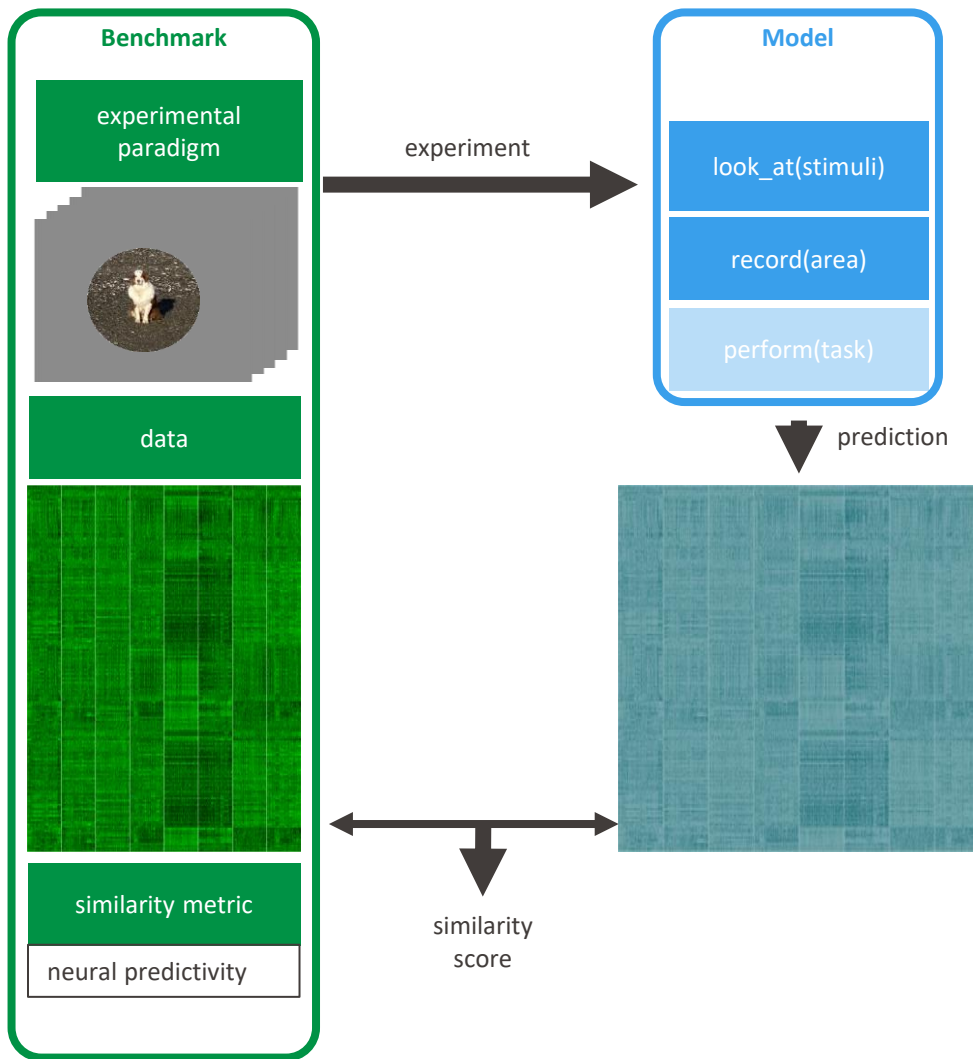
*Rajalingham\*, Issa\*, et al. (JNeuro 2018)*

EPFL

Neural data

video courtesy of Kailyn Schmidt, MIT

EPFL

Neural alignment = alignment between stimulus-matched recordings

**Neural benchmarks**

Brain recordings

Model units

stimuli

**Neural benchmarks**

neural predictivity

Brain recordings

Model units

stimuli

EPFL

**Neural benchmarks**

similarity metric

neural predictivity

Brain recordings

Model units

stimuli

fit

regression weights

correlation

predict held-out

*Yamins\*, Hong\*, et al. 2014*     *Schrimpf\*, Kubilius\*, et al. 2018*

# Model building

# Large scale architecture search and model comparison

model

primate

V1
V4
V2
IT

RGC   LGN   V1   V2   V4   IT   behavior

best mobilenet

best basenet

alexnet

Brain-Score

.55
.50
.45
.40
.35

r = 0.92

0   20   40   60   80
ImageNet performance (% top-1)

.56
.55
.54
.53
.52

cornet_s
densenet-169
resnet-101_v2
densenet-201
resnet-152_v2
xception
pnasnet_large
vgg-19
inception_v4
n.s.

70   72   74   76   78   80   82
ImageNet performance (% top-1)

Brain-Score

Schrimpf & Kubilius et al. 2018, 2020

# Brain-Score 100+ brain & behavior benchmarks, 300+ models
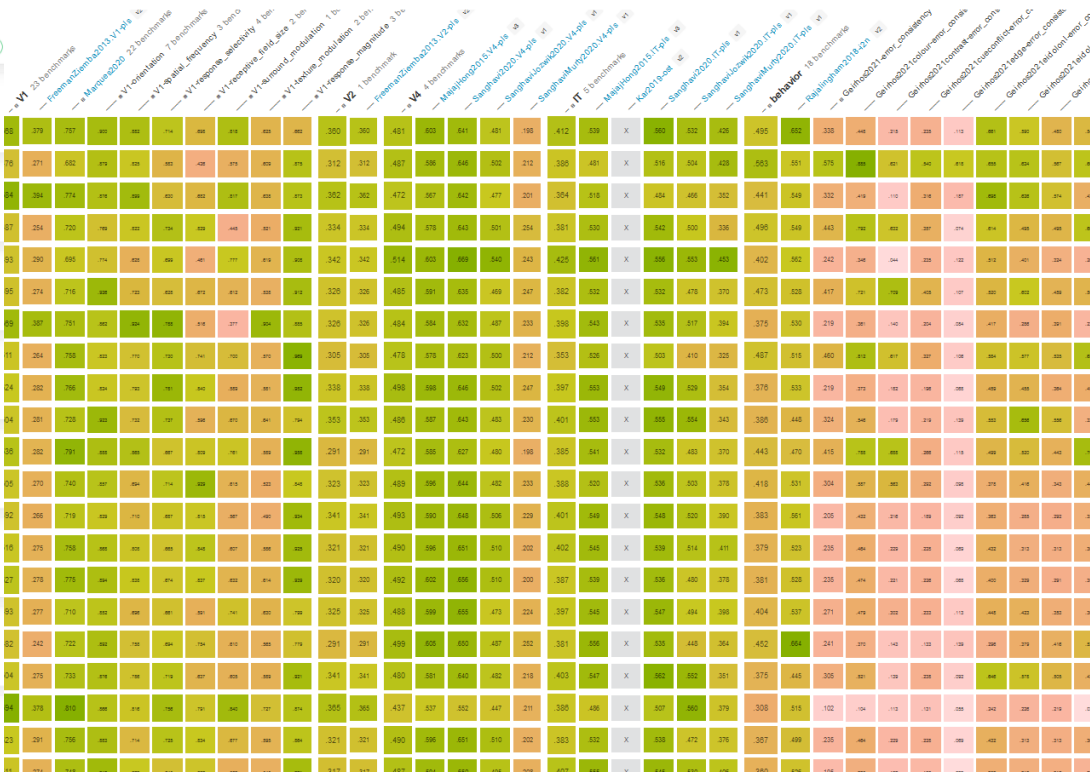
e.g. neural predictions for different image sets, distributional alignments such as spatial frequency, behavioral generalization, …
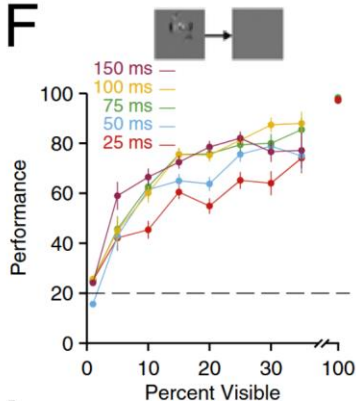
www.Brain-Score.org

# Recurrent processing in the visual system



F

150 ms —
100 ms —
75 ms —
50 ms —
25 ms —

G

150 ms —
100 ms —
75 ms —
50 ms —
25 ms —

Percent Visible

Percent Visible

Performance

Performance

Electrode in left fusiform gyrus (face-selective)

IFP (μV)

350
0
-350
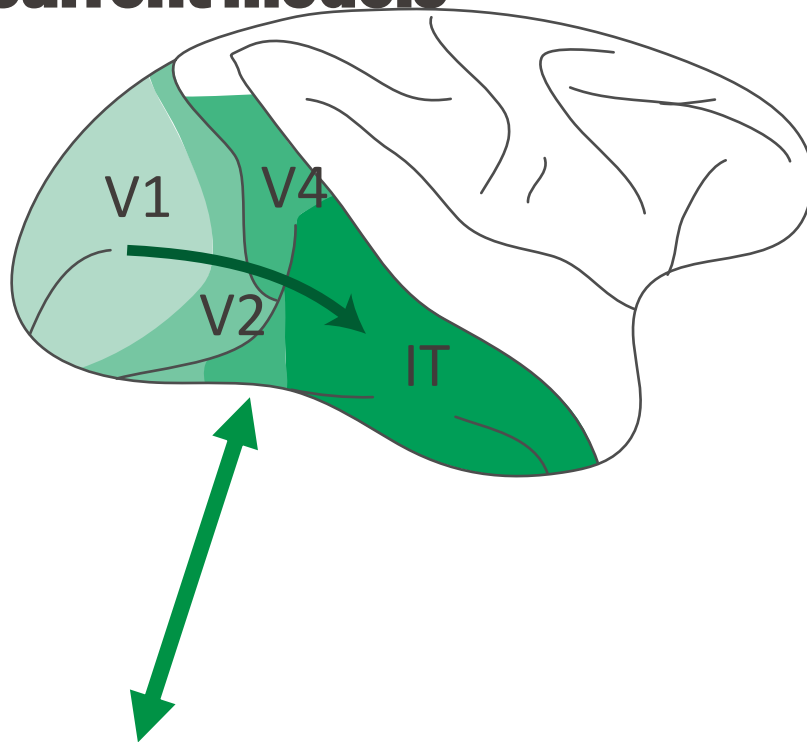
Time (ms)

Tang & Schrimpf & Lotter et al. 2018

- **Control images** are solved quickly
- **Challenge images** require more processing

"face"   "zebra"        "car"      "dog"    ...

behavioral accuracy

Decoding accuracy (d')

time from image onset

Kar et al. 2019

# Modeling recurrence: transform feed-forward networks into recurrent models

V1

V4

V2

IT

e.g. ResNet-101

$W_{1,1}$ $W_{1,2}$ $W_{1,3}$ $W_{1,...}$ $W_{1,5}$ $W_{1,6}$ $W_{2,1}$ $W_{2,2}$ $W_{2,3}$ $W_{2,...}$ $W_{2,11}$ $W_{2,11}$ $W_{3,1}$ $W_{3,2}$ $W_{3,3}$ $W_{3,...}$ $W_{3,31}$ $W_{3,32}$ $W_{4,1}$ $W_{4,2}$ $W_{4,3}$ $W_{4,...}$ $W_{4,11}$ $W_{4,12}$

# Transform feed-forward networks into recurrent models

e.g. He, Zhang, Ren, Sun (CVPR 2016)
Huang, Liu, van der Maaten, Weinberger (CVPR 2017)

Liang & Hu (CVPR 2015)    Liao & Poggio (arXiv 2016)    Tang*, Schrimpf*,
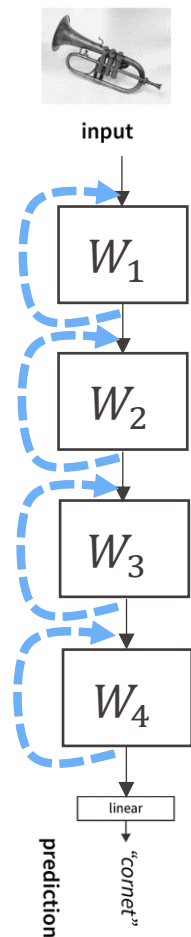Lotter* et al. (PNAS 2018)    Nayebi*, Bear*, Kubilius* et al. (NIPS 2018)

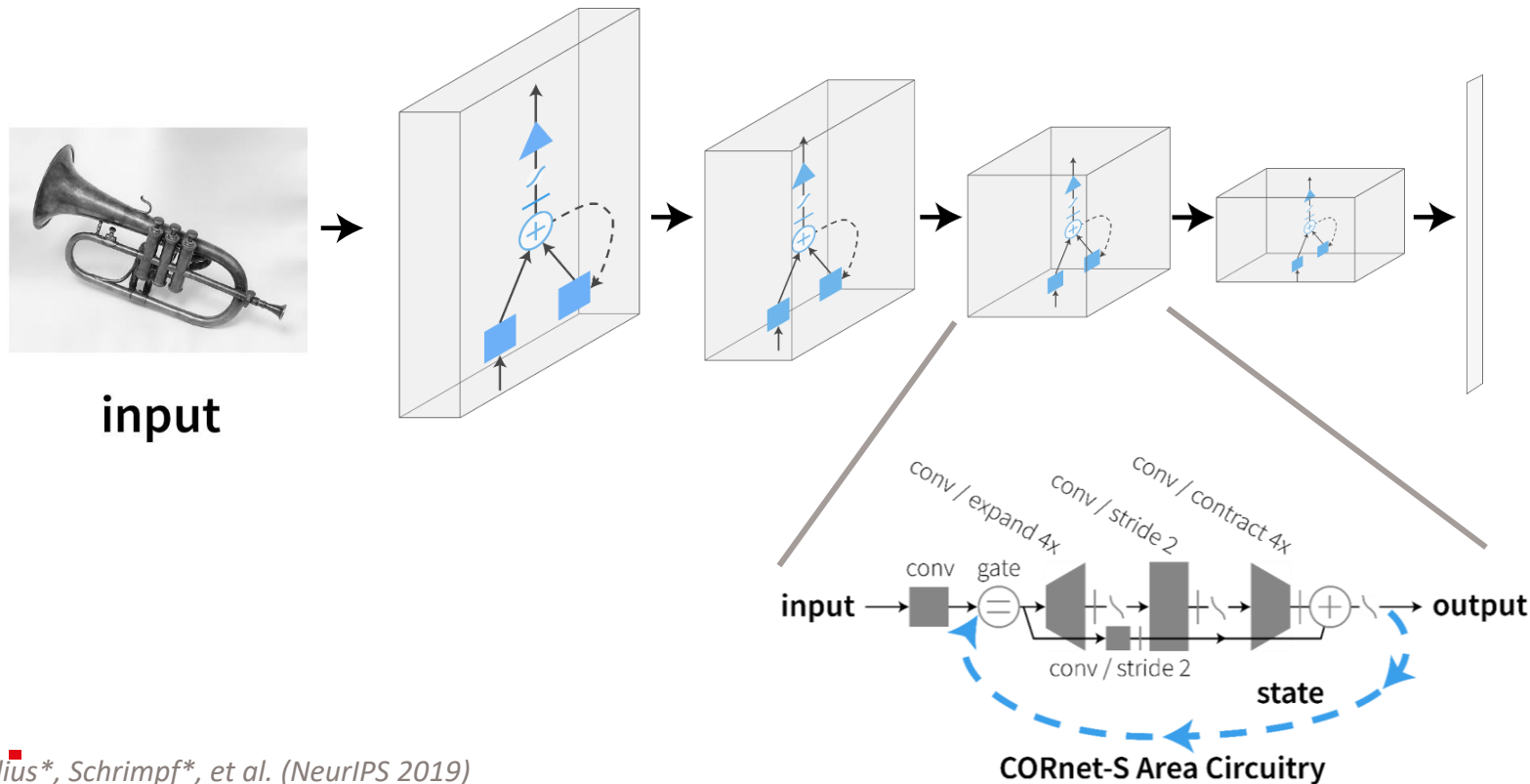# Modeling recurrence: transform feed-forward networks into recurrent models: CORnet



*Kubilius\*, Schrimpf\*, et al. (NeurIPS 2019)*

# Recurrent CORnet model: compact architecture via recurrence



CORnet-S Area Circuitry

*Kubilius\*, Schrimpf\*, et al. (NeurIPS 2019)*

# Recurrent model predicts temporal dynamics in IT

"face"   "zebra"   "car"   "dog"   …

behavioral accuracy

Decoding accuracy (d')

time from image onset

$IT_{monkey}$ object solution times

score: 0.3

$IT_{COR}$ object solution times

- Unlike feedforward models, CORnet-S can **predict neural responses over time**.

- i.e., when the **brain's IT is fast** to process images, **CORnet's IT-layer is also fast**

*Kar, Kubilius, Schmidt, Issa, DiCarlo (Nature Neuroscience 2019)*

*Kubilius\*, Schrimpf\*, et al. (NeurIPS 2019)*
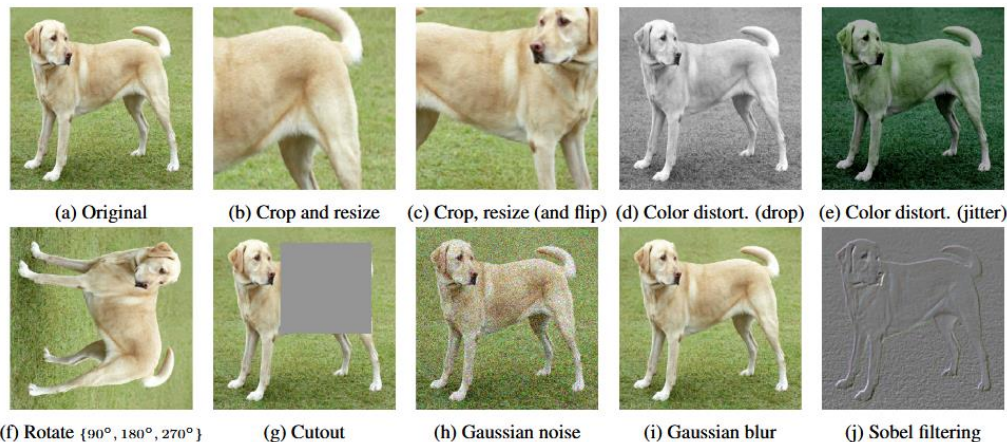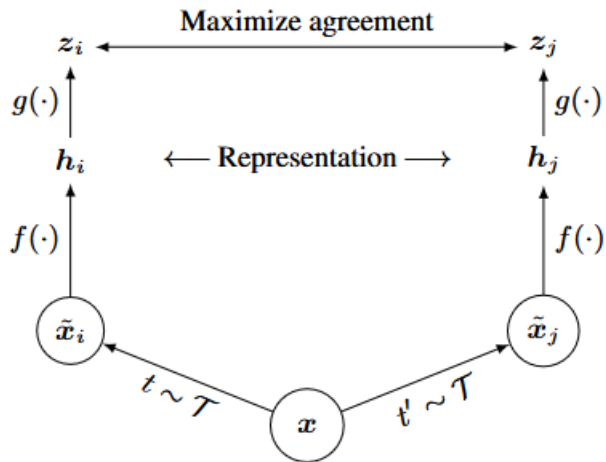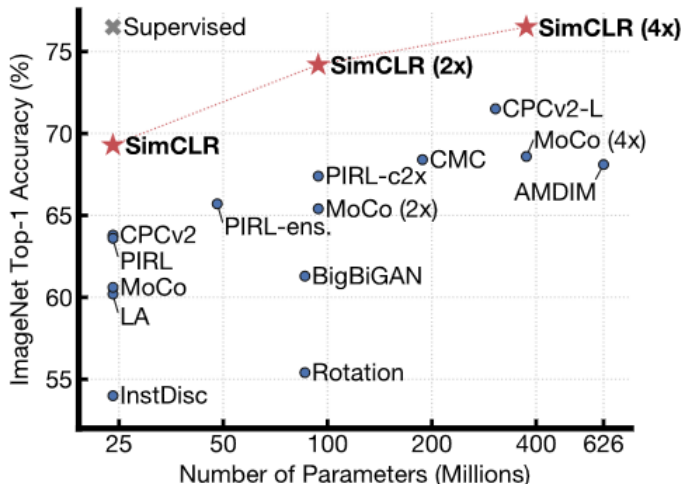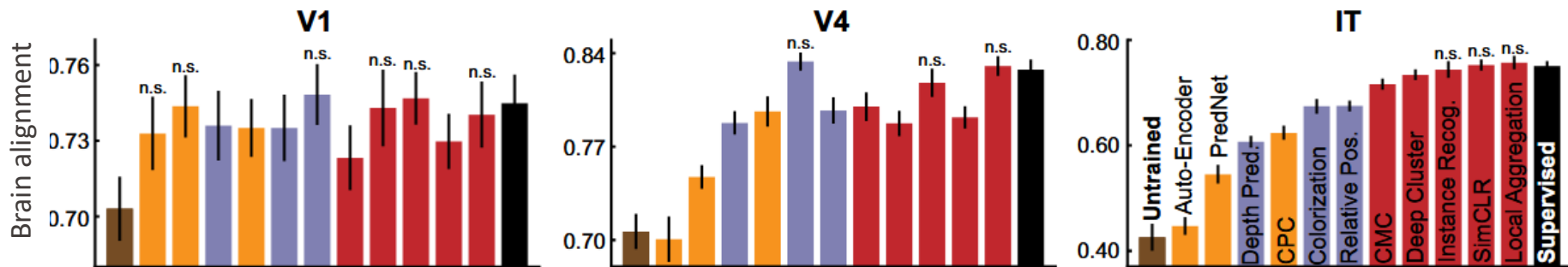
# Unsupervised learning with a contrastive loss

Figure 2. A simple framework for contrastive learning of visual representations. Two separate data augmentation operators are sampled from the same family of augmentations ($t \sim \mathcal{T}$ and $t' \sim \mathcal{T}$) and applied to each data example to obtain two correlated views. A base encoder network $f(\cdot)$ and a projection head $g(\cdot)$ are trained to maximize agreement using a contrastive loss. After training is completed, we throw away the projection head $g(\cdot)$ and use encoder $f(\cdot)$ and representation $h$ for downstream tasks.
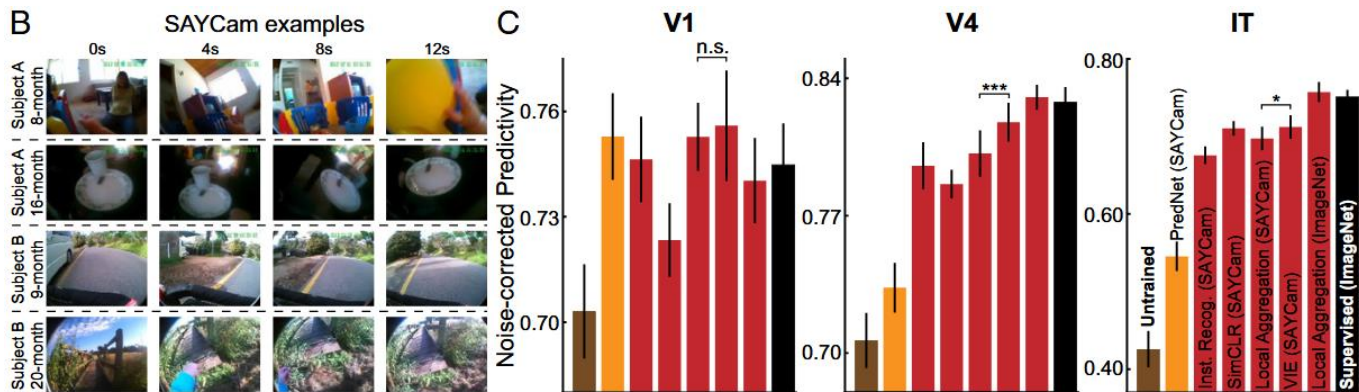
- Unsupervised approaches such as SimCLR encourage similar representations for similar inputs

- Performance rivals supervised learning
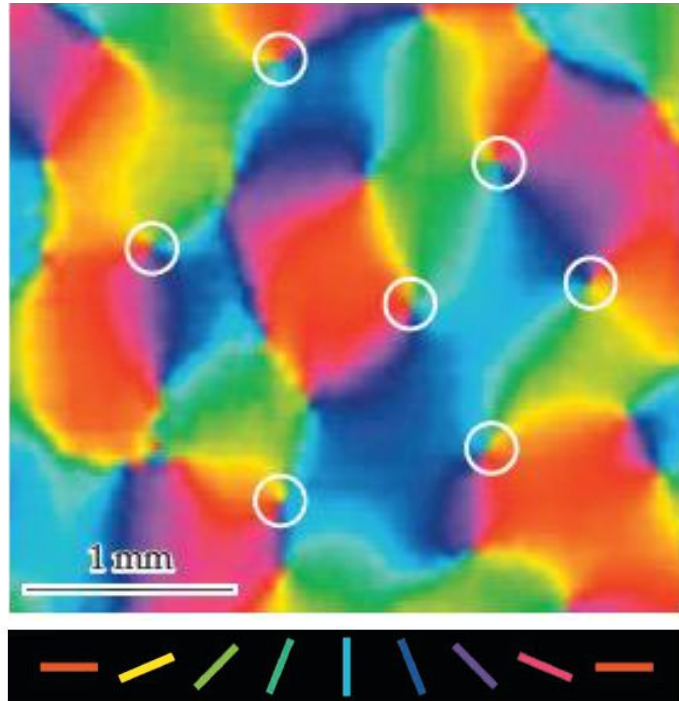
- <u>Chen et al. 2020</u>

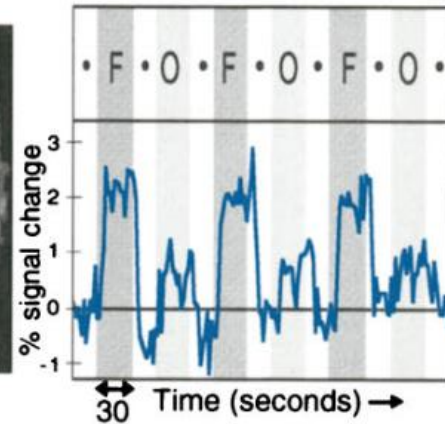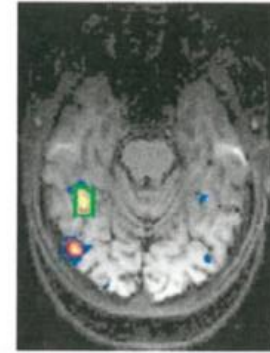# Unsupervised models also explain visual brain activity



- When trained on regular computer vision datasets (top) or developmental data streams (below, SAYCam), unsupervised models develop brain-like visual representations
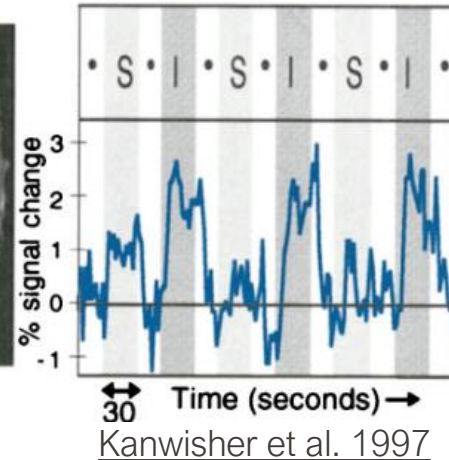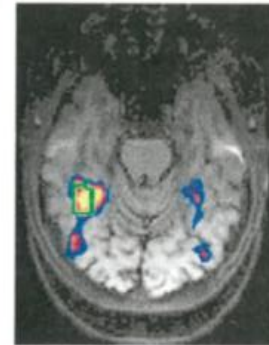


Zhuang et al. 2021

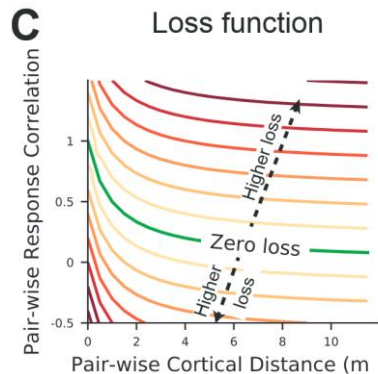# Neurons in cortex are topographically organized

**3a. Faces > Objects**

**3b. Intact Faces > Scrambled Faces**

Purves, chapter 11

Kanwisher et al. 1997

- A spatial loss term leads to brain-like clusters along the visual ventral stream (V1 to IT)



**A** Monkey IT

**C** Loss function

Macaque V1    TDANN

200μm

**F** Human VTC    TDANN

Faces
Bodies Places
Characters Objects

Lee et al. 2020

Margalit et al. 2024

# Topographic models enable the modeling of causal interventions, e.g. micro-stimulation

# Topographic models enable the modeling of causal interventions, e.g. micro-stimulation



stimulate near face-selective cluster

*Schrimpf et al. 2024*

# Topographic models enable the modeling of causal interventions, e.g. micro-stimulation



MODEL PREDICTION

BIOLOGICAL DATA (Afraz et al. 2006)

*Schrimpf et al. 2024*

# Synthesis of stimuli for neural population control

- Idea: model is fully differentiable, so we can set a desired target neural activity and update pixels in a way that they elicit the target state (according to model predictions)

# Model-guided synthesis non-invasively controls neural activity

EPFL

- This procedure works!
- We can generate stimuli that drive neural activity beyond the typical range
- This is a non-invasive control procedure

Bashivan*, Kar*, DiCarlo 2019

# Generating "exciting" stimuli without a pre-trained encoding model

Walker et al. 2019

Ponce & Xiao & Schade et al. 2019

# Adversarial attacks in computer vision

- Models are fooled by small, imperceptible perturbations (white box adversarial attacks)

- Protection technique: train on adversarial images "adversarial training" (very costly)

Graffiti

Egyptian Cat
*Milou*

Perturbation
$(\|\delta\|_2 = 0.60)$

Sleeping Bag

Image from *Dapello*, Marques*, et al. (NeurIPS 2021)*

*Szegedy et al. (ICLR 2014)*
*Eykholt*, Evtimov*, et al. (CVPR 2018)*

# Adversarial attacks on the brain

- Prevalent view: only computational models are susceptible to adversarial attacks

- But: can synthesize images that also fool IT neurons

Guo et al. 2022

# Adversarial attacks on behavior

- Using a robustified model (**trained with adversarial attacks**), can change images in a way that change the decision of humans

**Vanilla guide model**

**Robustified guide model**

Normalized behavioral disruption [% errors]

100

0

Vanilla model

Gap

Humans

Unperturbed

Robustified model

Gap closure

Humans

Unperturbed

$10^{-1}$  $10^0$  $10^1$

$10^{-1}$  $10^0$  $10^1$

Perturbation pixel budget ($\epsilon$) [$\ell_2$-norm]

Gaziv*, Lee*, DiCarlo 2023

# Adversarial attacks on behavior

**start images**

**Example target categories**
'insect'    'primate'    'fish'



Gaziv*, Lee*, DiCarlo 2023

# Metamers

- Model metamers: "stimuli whose activations within a model stage are matched to those of a natural stimulus"



Feather et al. 2023

**Metamers of standard models are not recognizable by humans**

Feather et al. 2023

# Adversarial training makes metamers human-recognizable

**b** ResNet50, adversarial training

Proportion correct (human)

*N* = 20

Chance performance (1/16)

**c** AlexNet, adversarial training

*N* = 20

$L_2$ ($\varepsilon$ = 3) adversarial perturbations
$L_\infty$ ($\varepsilon$ = 4/255) adversarial perturbations
$L_\infty$ ($\varepsilon$ = 8/255) adversarial perturbations

Standard training
$L_2$ ($\varepsilon$ = 3) random perturbations
$L_\infty$ ($\varepsilon$ = 8/255) random perturbations

Feather et al. 2023

Example metamers (ResNet50, $L_2$ perturbations)

**d**

natural image | conv1_relu1 | layer1 | layer2 | layer3 | layer4 | avgpool | final

Standard training

$L_2$ ($\varepsilon$ = 3) Adversarial perturbations

# Models of auditory processing

- Jointly optimize CNN for word + genre recognition tasks



**A** Word recognition task
Excerpted speech + Background noise
587-way AFC: Which word (at 1 sec.)?
2 sec.

Musical genre task
Excerpted music + Background noise
41-way AFC: Which genre?
2 sec.

**E** Best-performing deep neural network
conv1 norm1 pool1 conv2 norm2 pool2 conv3 conv4 conv5 pool5 fc6 fc_top
Word classifier
Genre classifier
Example first-layer filters

**F** Baseline model: Spectrotemporal filter bank
Example spectrotemporal filters

Kell & Yamins et al. 2018

# Task-optimized model exhibits human-like behavior



EPFL

**A** Human word recognition

**B** Network word recognition

**C** Word recognition: Human vs. network

Background type: □ Music □ Auditory scene □ Speaker-shaped noise □ 2-speaker babble □ 8-speaker babble

**G** Genre recognition: Humans

**H** Genre recognition: Network

**I** Genre recognition: Human & Network

Background type: □ Auditory scene □ Music-shaped noise □ 2-speaker babble □ 8-speaker babble

- Model closely predicts human performance patterns, especially for word recognition tasks

- Less behaviorally-aligned for genre recognition

Kell & Yamins et al. 2018

# Task-optimized audio model predicts fMRI responses

- The task-optimized audio model predicts brain activity in auditory cortex better than baseline models

**Trained network** (selected architecture, trained filters)
**Spectrotemporal model**
**Random-filter network** (selected architecture, untrained filters)
**Random-filter network** (unselected architectures, untrained filters)
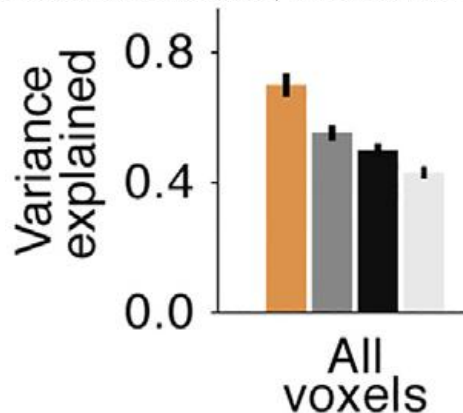
Variance explained: 0.8, 0.4, 0.0

**All voxels**

**A**

**165 everyday sounds:**

person screaming
velcro
whistling
frying pan sizzling
alarm clock
cat purring
guitar riff
... etc. ...

$v_i$: single voxel's response to all 165 natural sounds

sounds

$v_i = Fw$
Regularized linear regression

**Region of interest**

F: features from one network layer

sounds

units →

Kell & Yamins et al. 2018

# Model Predicts Hierarchical Organization in Human Auditory Cortex



- Black outline: sub-divisions of primary auditory cortex

- Primary auditory cortex best explained by earlier layers

- Later layers best explain non-primary areas

Kell & Yamins et al. 2018

# Take-home messages

- Unsupervised training yields brain-like representations

- Including a spatial loss term leads to topographic models that reproduce the spatio-functional organization in the brain. These models can predict the behavioral effects of neural interventions

- Encoding models of brain function allow for the synthesis of images to control neural activity. Can create images akin to adversarial attacks

- Yet, models often view things as identical that appear very different to humans (metamers)

- Very similar ideas from vision work for models of auditory processing