# Invariant & hierarchical computation in human auditory cortex

## Alex Kell

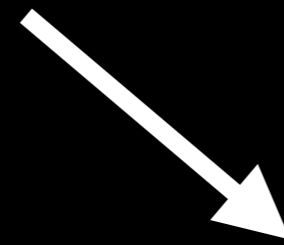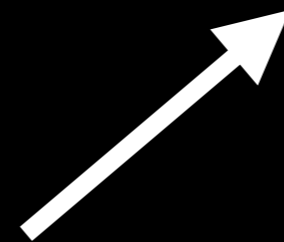**2018.07.17 :: CSGF Program Review**

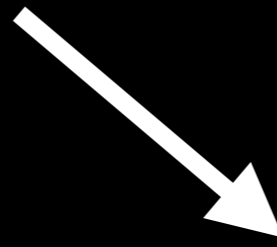time ➝



time ➝

2

time ⟶

time ⟶

What was said?
Who said it?
How did they feel when they said it?
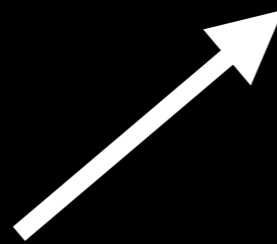
What caused the sound?
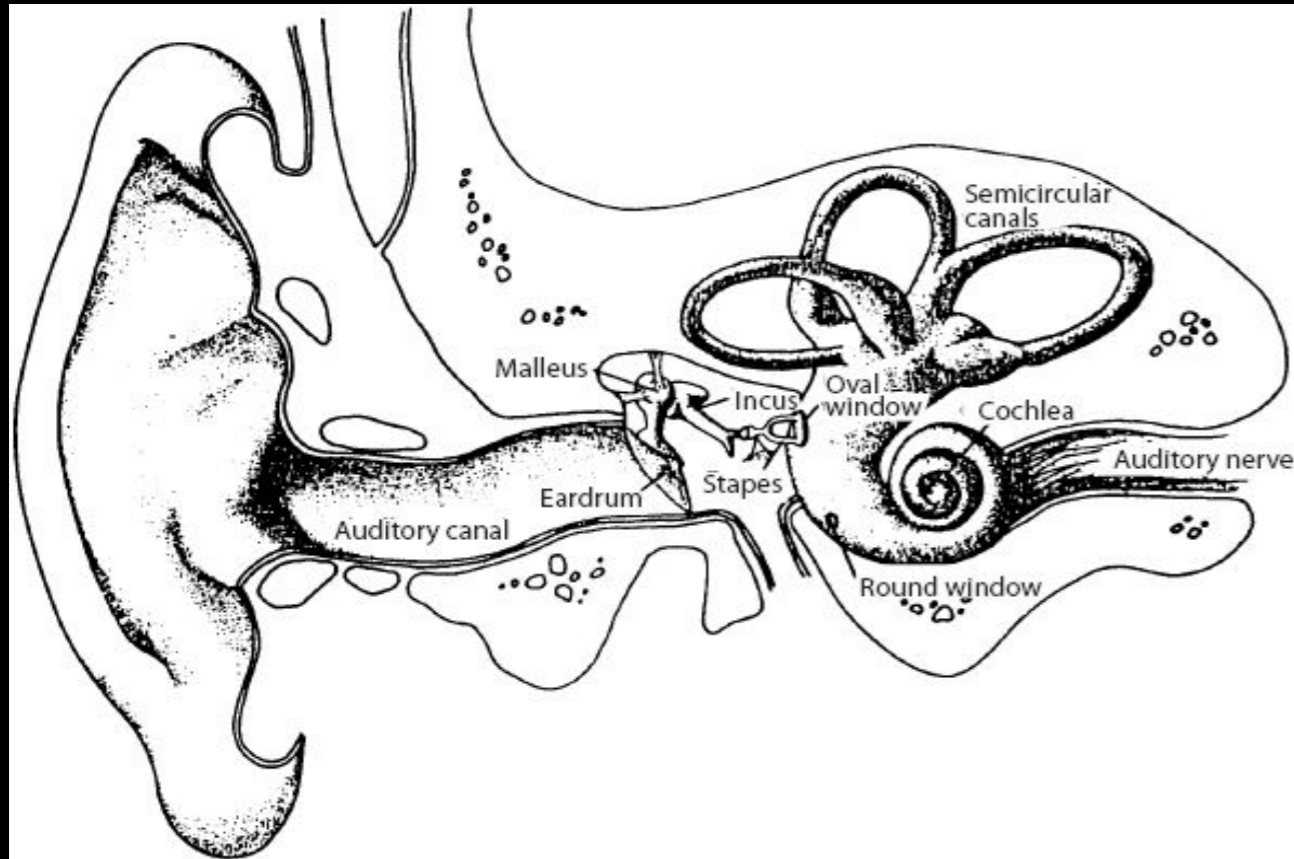Where?

What was said?
Who said it?
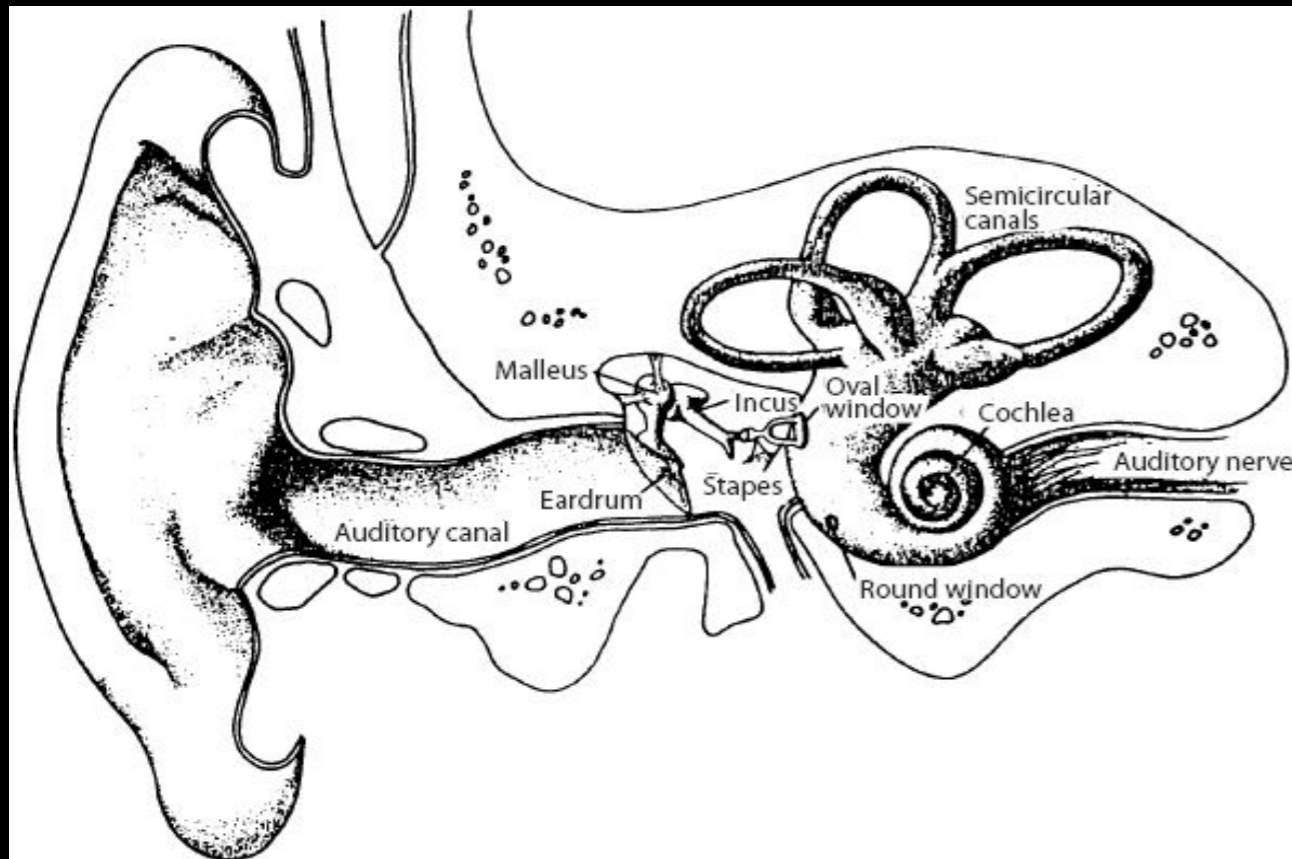How did they feel when they said it?

What caused the sound?
Where?

**How does the brain extract behaviorally relevant information from these waveforms?**

# Peripheral auditory system: well characterized...

# Peripheral auditory system: well characterized...



**... but auditory cortex is poorly understood.**
**(Particularly in humans.)**

# TODAY:

**Basic questions about functional organization of human auditory cortex.**

# TODAY:

## Is there a hierarchical organization?

# TODAY:

Is there a hierarchical organization?

If so, how many stages?

# TODAY:

Is there a hierarchical organization?

If so, how many stages?

What do different stages do?

# TODAY:

Is there a hierarchical organization?

If so, how many stages?

What do different stages do?

Use modeling to generate specific hypotheses in a principled manner

# A Task-Optimized Neural Network Replicates Human Auditory Behavior, Predicts Brain Responses, and Reveals a Cortical Processing Hierarchy
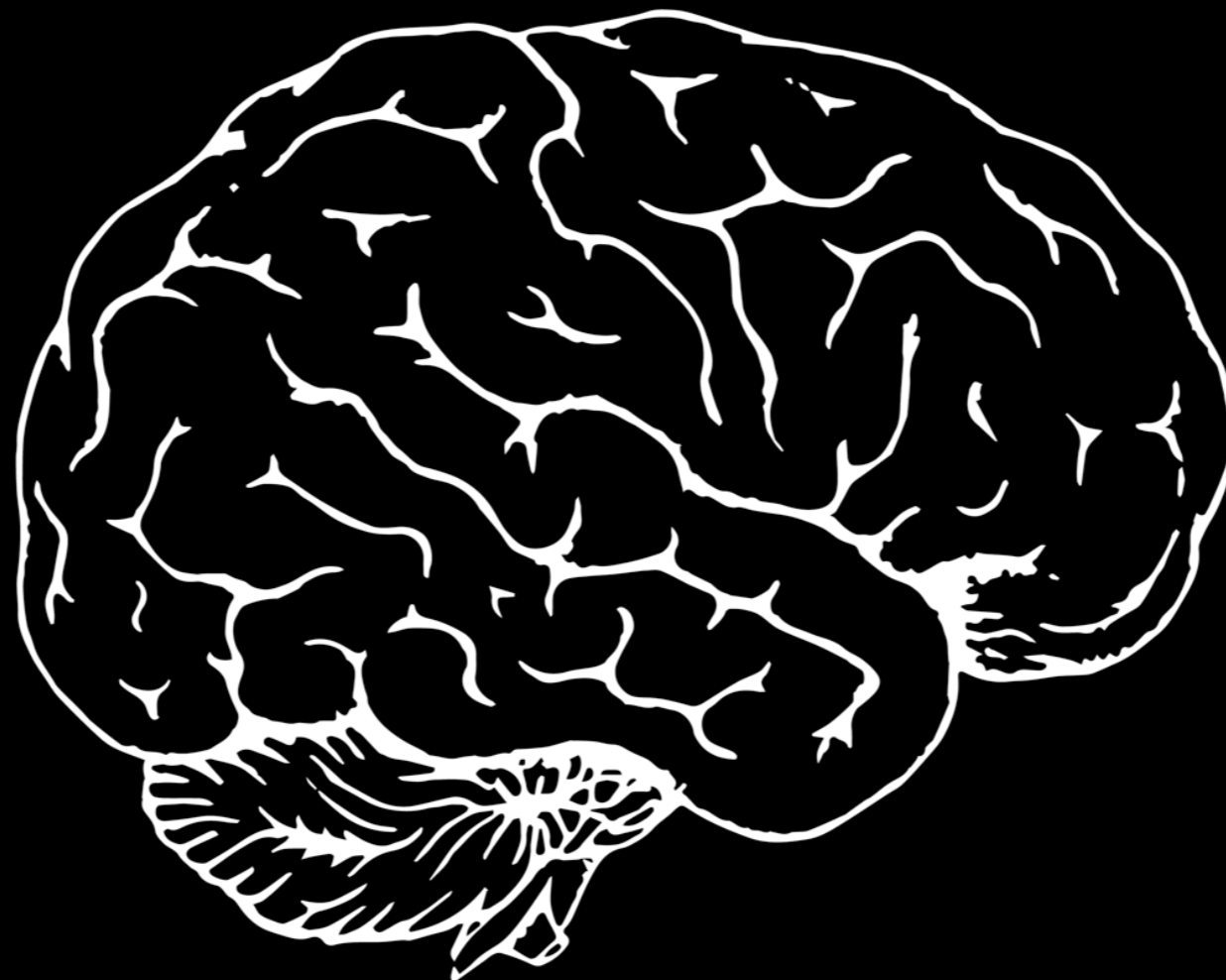
## Highlights

- A deep neural network optimized for speech and music tasks performed as well as human listeners

- The optimization produced separate music and speech pathways after a shared front end

## Authors

Alexander J.E. Kell, Daniel L.K. Yamins,
Erica N. Shook,
Sam V. Norman-Haignere,
Josh H. McDermott

**Work with:**
**Dan Yamins, Erica Shook, Sam Norman-Haignere,**
**and Josh McDermott**

# How to build better models of auditory cortex?

# How to build better models of auditory cortex?

time →

What was said?
Who said it?
How did they feel when they said it?

What caused the sound?
Where?

# How to build better models of auditory cortex?



time →

MODEL

What was said?
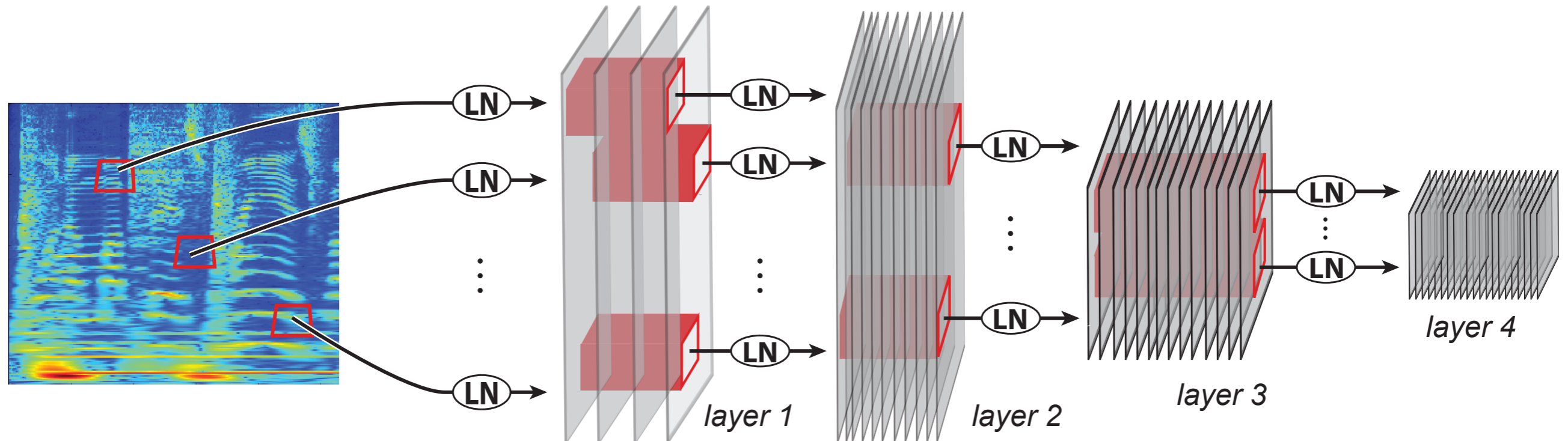Who said it?
How did they feel when they said it?

What caused the sound?
Where?

# Recent machine learning advances: Deep learning

# Recent machine learning advances: Deep learning



# Hierarchical convolutional neural networks (CNNs)
(Fukushima, 1980; Lecun et al., 1989; Krizhevsky et al., 2012; Yamins, Hong, et al., 2014; etc.)

# KEY HYPOTHESIS:

**A model optimized to perform real-world auditory tasks may converge to brain-like computations**

## KEY HYPOTHESIS:

**A model optimized to perform real-world auditory tasks may converge to brain-like computations**

**Approach pioneered in the visual cortex**
(Yamins, Hong, et al. 2014; Cadieu et al. 2014; Hong, Yamins, et al. 2016)

## KEY HYPOTHESIS:

# A model optimized to perform real-world auditory tasks may converge to brain-like computations

## Approach pioneered in the visual cortex
(Yamins, Hong, et al. 2014; Cadieu et al. 2014; Hong, Yamins, et al. 2016)

## Potentially:
## Particularly useful in auditory cortex

# Unsatisfying aspects of deep learning as a neuroscience model

# Unsatisfying aspects of deep learning as a neuroscience model

- Unrealistic amount of (supervised) training data.

- Unrealistic learning rule (backprop).

- Discriminative models (rather than generative).

- etc.

# Unsatisfying aspects of deep learning as a neuroscience model

- **Unrealistic amount of (supervised) training data.**

- Unrealistic learning rule (backprop).

- Discriminative models (rather than generative).

- etc.

# Network optimization:
# Real-world tasks

# Network optimization:
# Real-world tasks

... that have large labelled datasets.

# Network optimization:
# Real-world tasks

## ... that have large labelled datasets.

# Network optimization:
# Real-world tasks

## ... that have large labelled datasets.



Word recognition task

Excerpted speech + Background noise

587-way AFC:
Which word (at 1 sec.)?
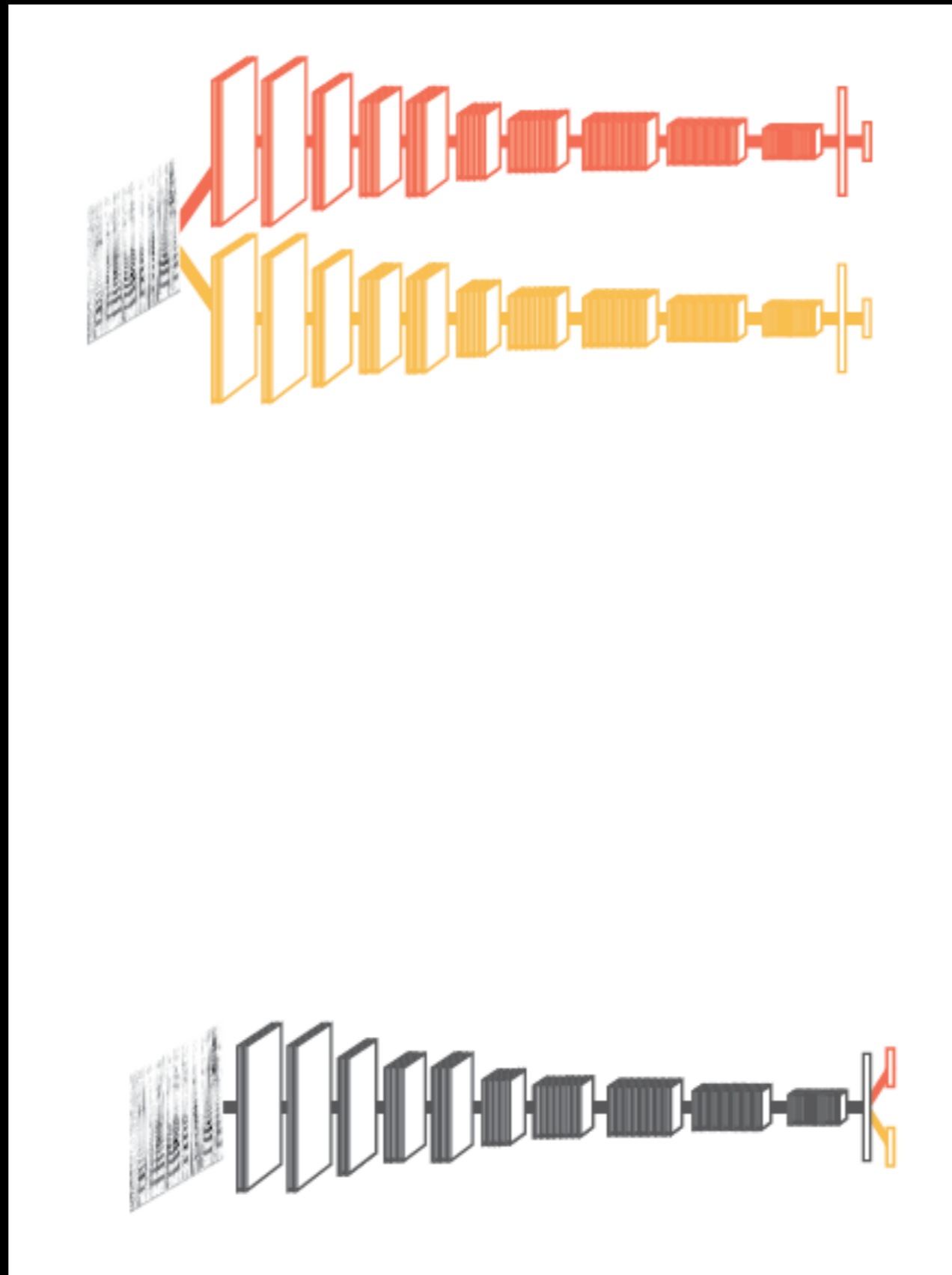
2 sec.

Musical genre task

Excerpted music + Background noise

41-way AFC:
Which genre?

2 sec.

# Network optimization: Architecture search

# Network optimization: Architecture search

# Network optimization: Architecture search
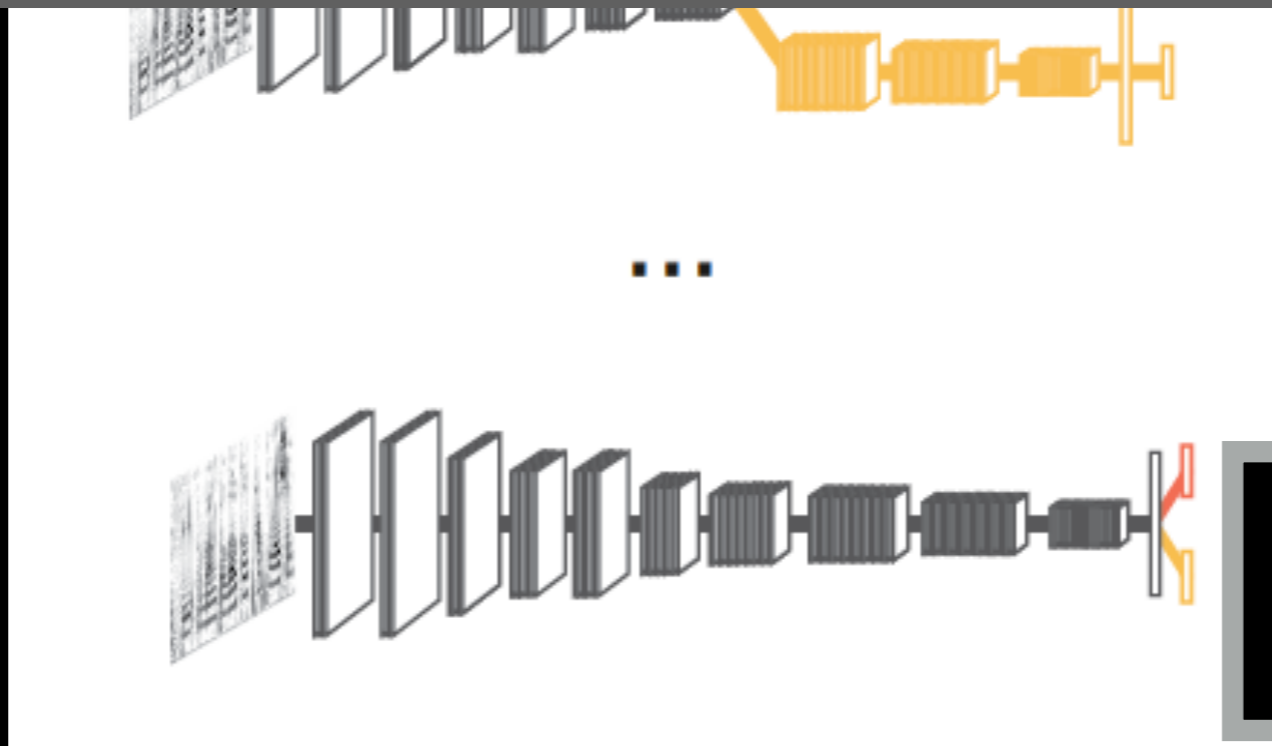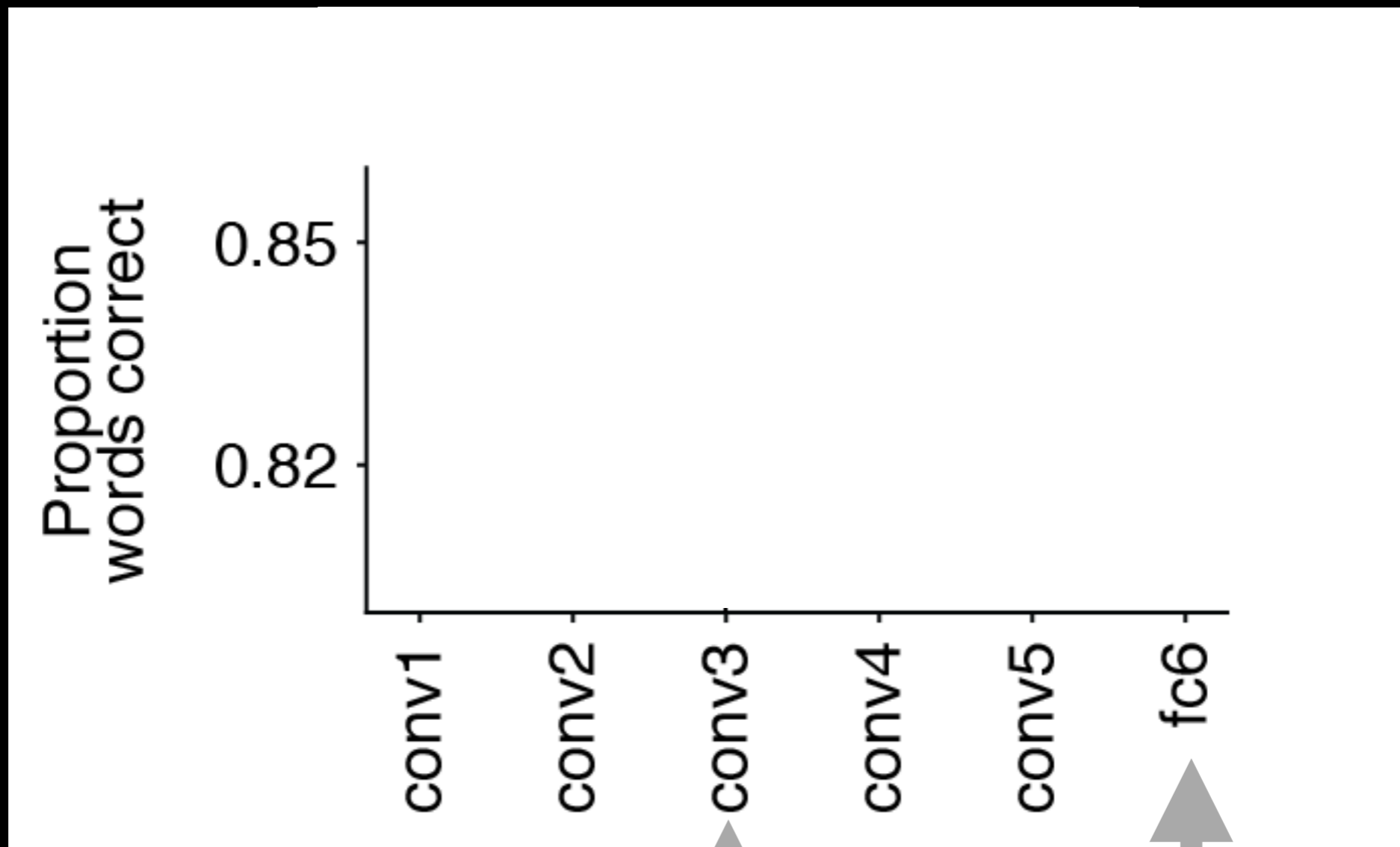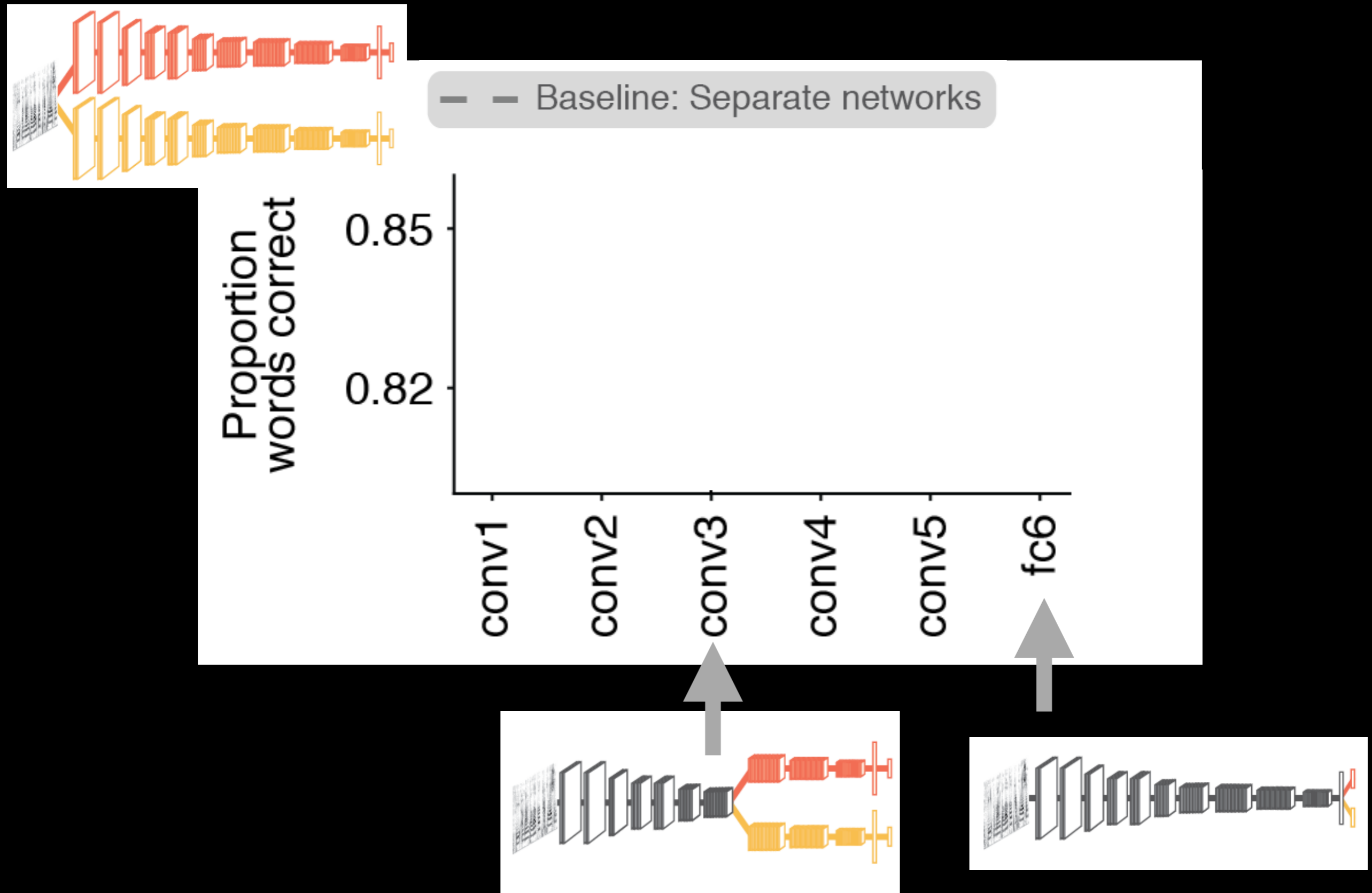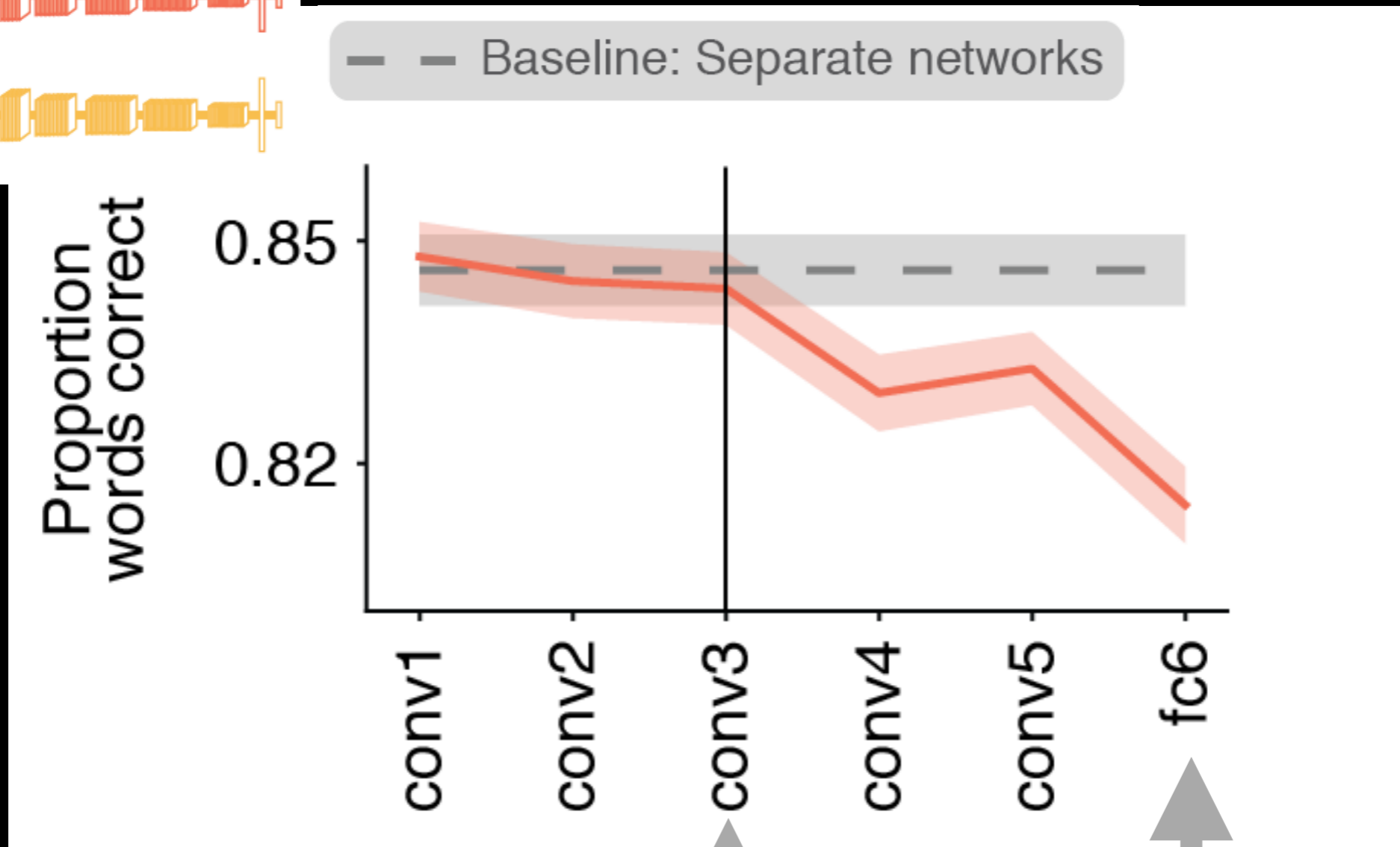
# Network optimization: Architecture search

# Network optimization: Architecture search



More parameters

Fewer parameters

# Network optimization: Architecture search



**More parameters**

**How many layers can be shared without a detriment in task performance?**

...

**Fewer parameters**

# Network optimization: Architecture search
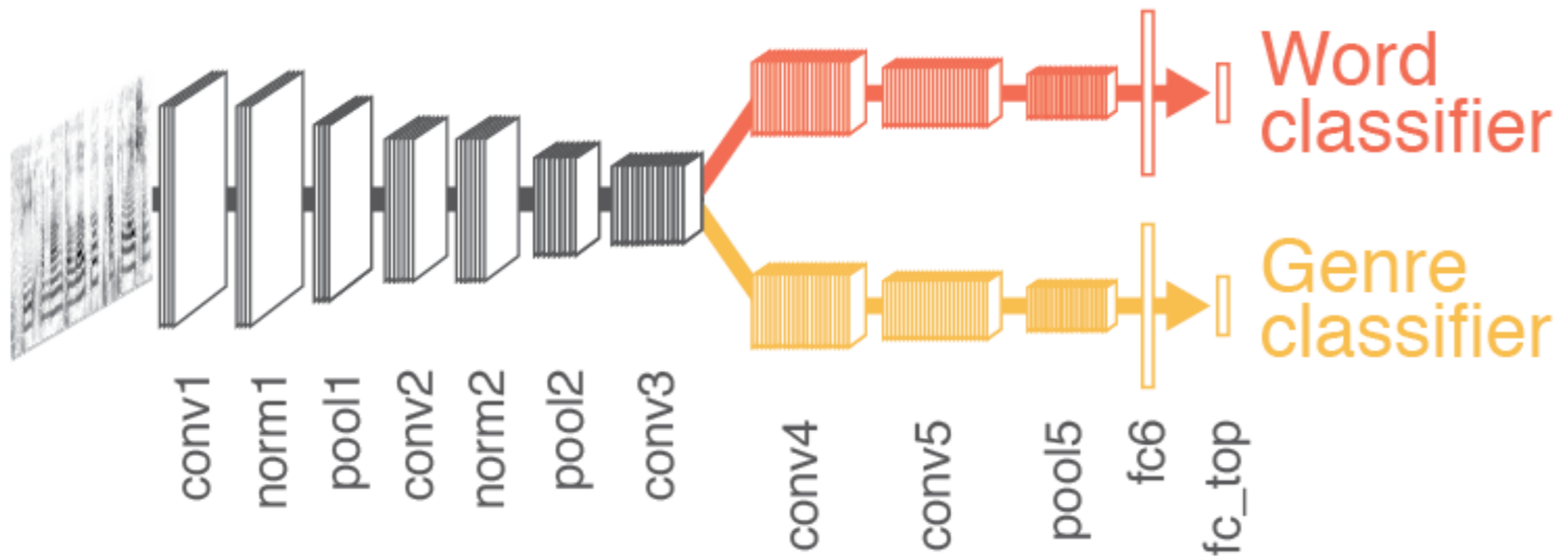
# Network optimization: Architecture search

# Network optimization: Architecture search
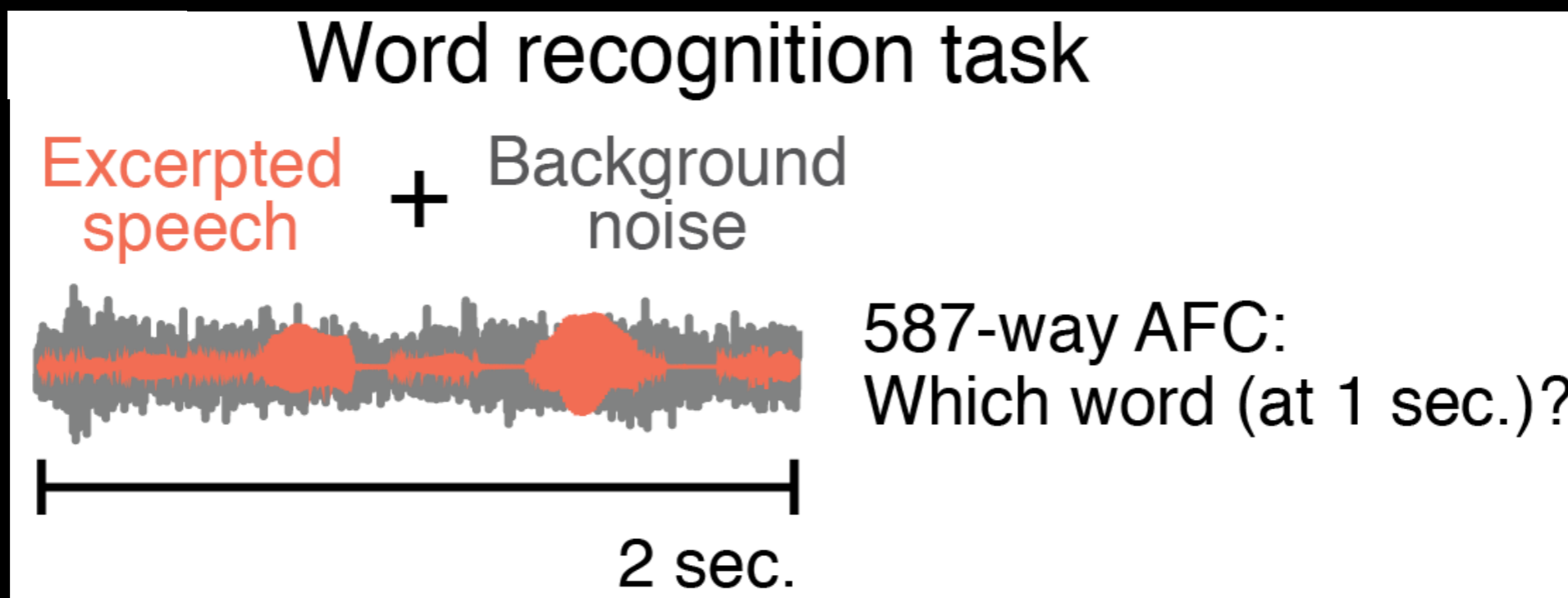
# Network optimization: Resulting network



Best-performing deep neural network

Word classifier

Genre classifier

conv1 · norm1 · pool1 · conv2 · norm2 · pool2 · conv3 · conv4 · conv5 · pool5 · fc6 · fc_top

Example first-layer filters

# Comparing human & model behavior

# Comparing human & model behavior



**26 conditions:**
**5 background types x 5 signal-to-noise ratios (SNRs)**
**+ noiseless**

# CNN & human psychophysics



Word psychophysics:
Humans

Proportion correct — 1.0, 0.5, 0.0
SNR (dB) — -9, -6, -3, 0, 3, Inf.

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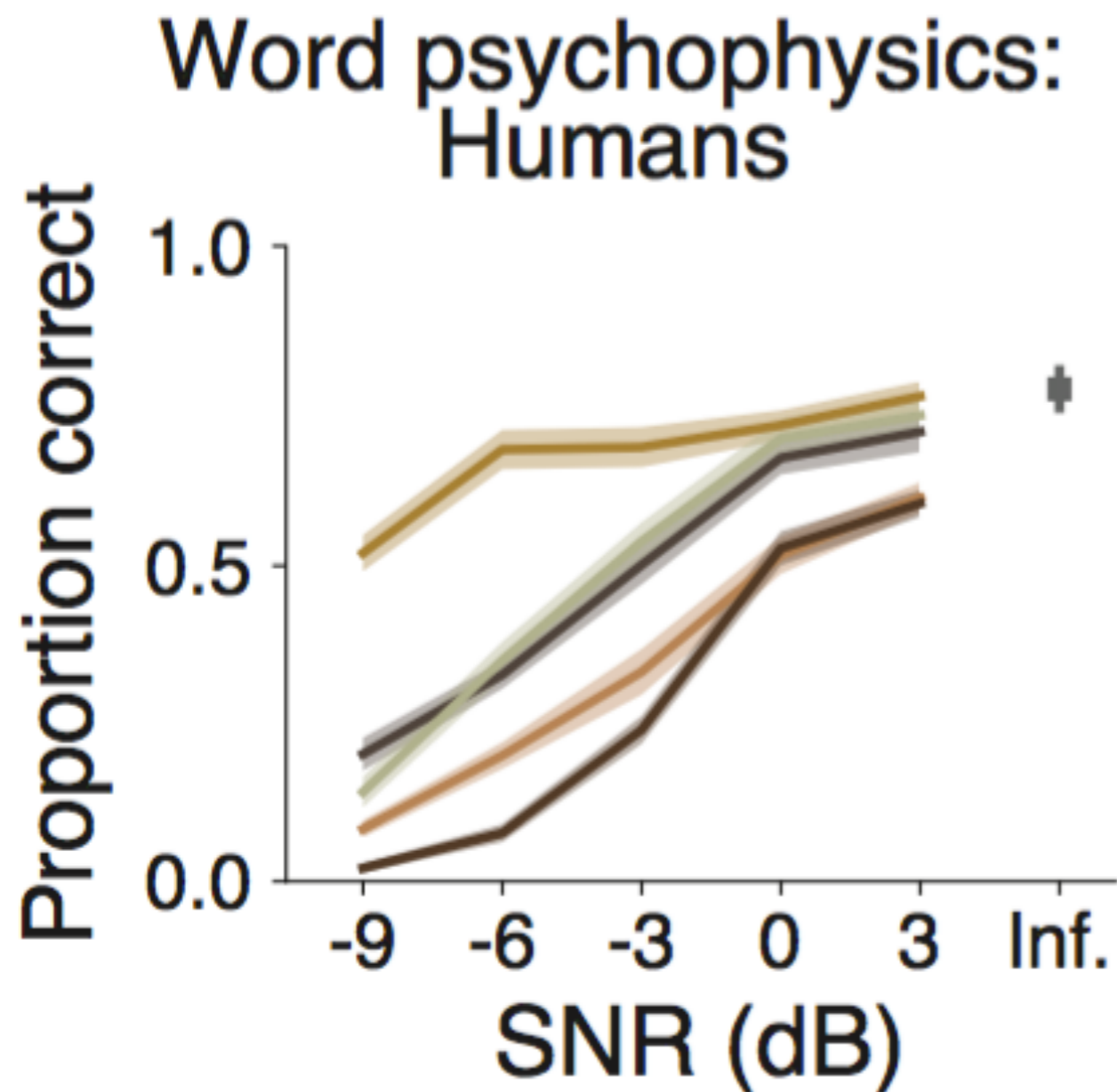- ☐ 8-speaker babble

# CNN & human psychophysics



Word psychophysics: Humans

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

# CNN & human psychophysics



Word psychophysics: Humans

Word psychophysics: Network

Proportion correct — SNR (dB): -9, -6, -3, 0, 3, Inf.

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

# CNN & human psychophysics

**NOTE:**

**CNN optimized ONLY for task performance
NOT optimized to behave similarly to humans**

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

# CNN & human psychophysics

**NOTE:**

CNN optimized ONLY for task performance
NOT optimized to behave similarly to humans

**POTENTIAL REASONS FOR SIMILARITY:**

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

# CNN & human psychophysics

## NOTE:

**CNN optimized ONLY for task performance**
**NOT optimized to behave similarly to humans**

## POTENTIAL REASONS FOR SIMILARITY:

1. Both network & humans near optimal?

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

# CNN & human psychophysics

## NOTE:

CNN optimized ONLY for task performance
NOT optimized to behave similarly to humans

## POTENTIAL REASONS FOR SIMILARITY:

1. Both network & humans near optimal?

2. Algorithmic similarities between net & humans?

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

# Using this model to predict cortical responses to natural sounds

# Using this model to predict cortical responses to natural sounds

## Measure fMRI responses to 165 natural sounds*

| | | |
|---|---|---|
| person screaming | road traffic | guitar |
| man speaking | zipper | coughing |
| flushing toilet | cellphone vibrating | crumpling paper |
| pouring liquid | water dripping | siren |
| tooth-brushing | scratching | splashing water |
| woman speaking | car windows | computer speech |
| car accelerating | telephone ringing | alarm clock |
| biting and chewing | chopping food | walking with heels |
| laughing | telephone dialing | vacuum |
| typing | girl speaking | wind |
| car engine starting | car horn | boy speaking |
| running water | writing | chair rolling |
| breathing | computer startup sound | rock song |
| keys jangling | background speech | door knocking |
| dishes clanking | songbird | dog barking |
| … | … | … |

*Norman-Haignere, Kanwisher, McDermott <u>Neuron</u> 2015 (Thanks!)

165 everyday sounds:

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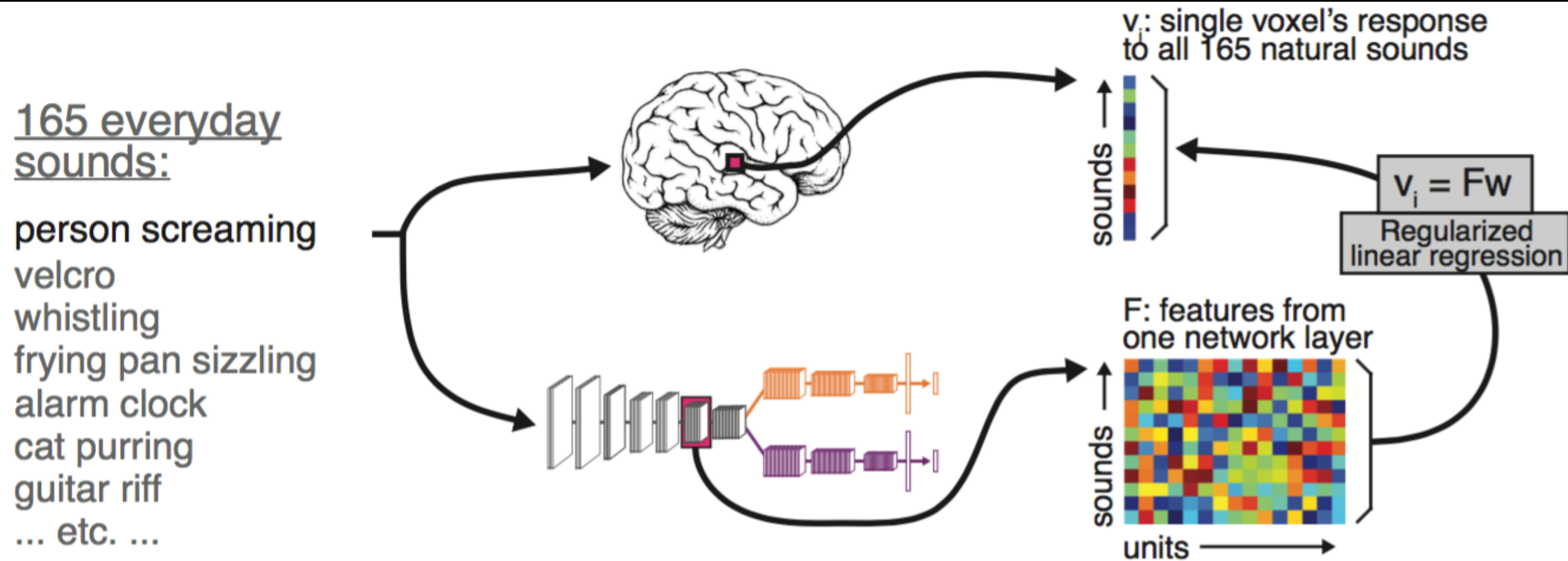frying pan sizzling
alarm clock
cat purring
guitar riff
... etc. ...

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

sounds

**Each voxel:**
**Mean response to each of 165 sounds**

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

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

F: features from one network layer

sounds

units →

# CNN as encoding model

Each voxel = weighted sum of units in a given layer

# CNN as encoding model

## Each voxel = weighted sum of units in a given layer



Cross-validated regularized linear regression
to predict voxel's response

# CNN as encoding model

## Each voxel = weighted sum of units in a given layer



**Cross-validated regularized linear regression
to predict voxel's response**

**Dependent measure: Variance explained**

# CNN as encoding model

## Each voxel = weighted sum of units in a given layer



Cross-validated regularized linear regression
to predict voxel's response

Dependent measure: Variance explained

Baseline:
Identical procedure with a spectrotemporal filter model

# Variance explained across all of auditory cortex

# Variance explained across all of auditory cortex

# Variance explained across all of auditory cortex

# Variance explained across all of auditory cortex

# Organization of human auditory cortex outside of primary areas?

# Organization of human auditory cortex outside of primary areas?

## A proposal from macaque anatomy:
### Tripartite hierarchical organization



Tramo et al. (1999)

**Evidence mostly anatomical**

# A measure of hierarchy?

# CNN architecture:
# Hierarchical and feedforward

# CNN architecture:
# Hierarchical and feedforward



Which layer best predicts each voxel's response?
A measure of "complexity"

Best-predicting network layer for each voxel

Best-predicting network layer for each voxel

Layer: ■ conv3 or lower ■ conv4 ■ conv5 or higher

Best-predicting network layer for each voxel
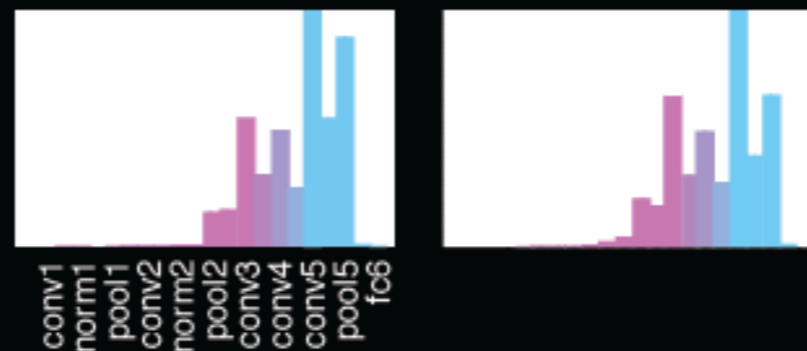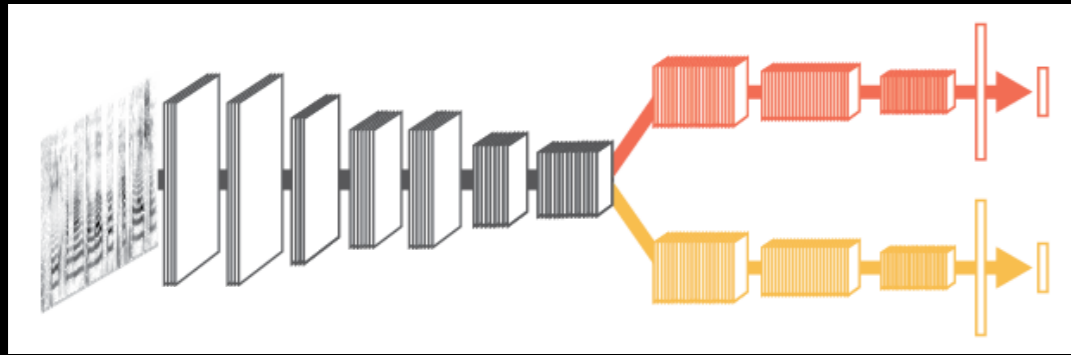
Layer: ■ conv3 or lower  ■ conv4  ■ conv5 or higher

RH    LH

70

Best-predicting network layer for each voxel

Layer: conv3 or lower — conv4 — conv5 or higher

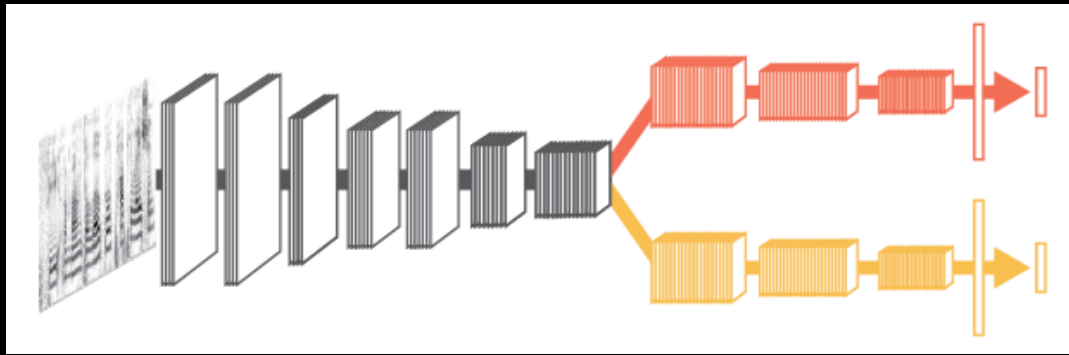**Network reveals hierarchical organization in human auditory cortex**
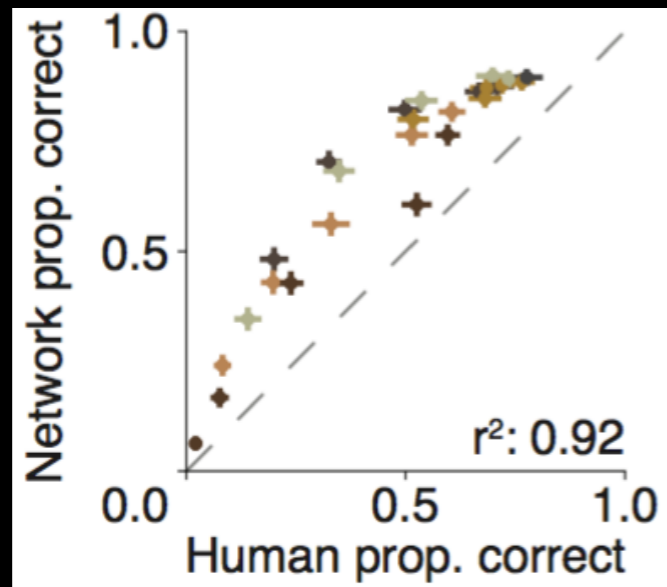
RH          LH
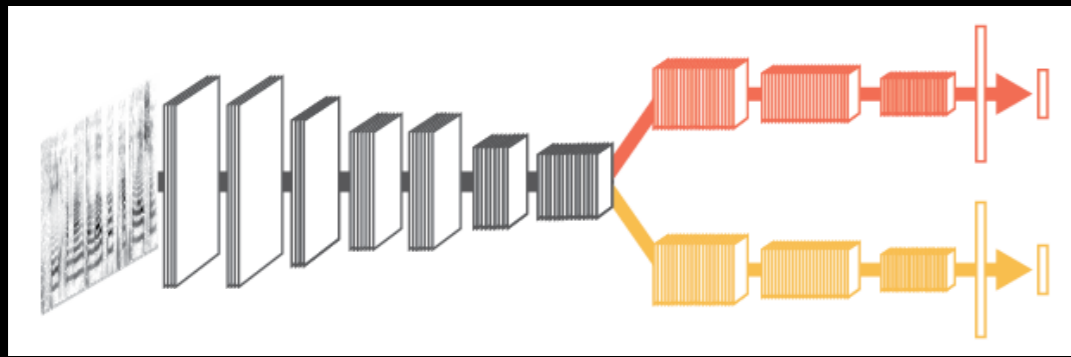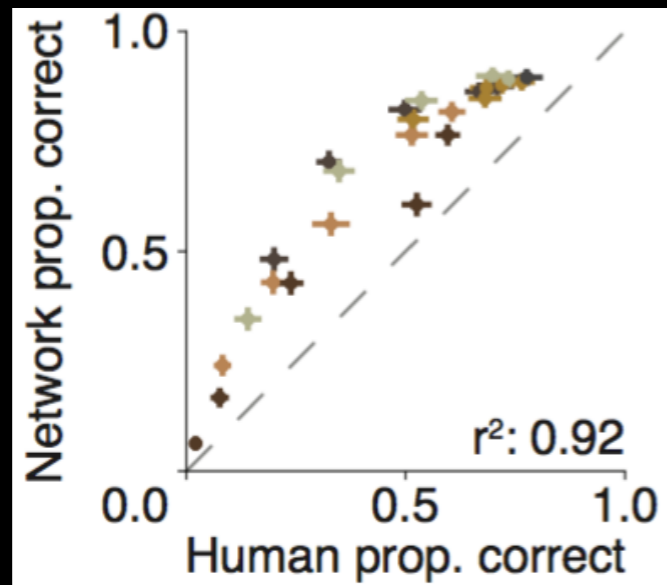
**Introduced multi-task networks as neural models**

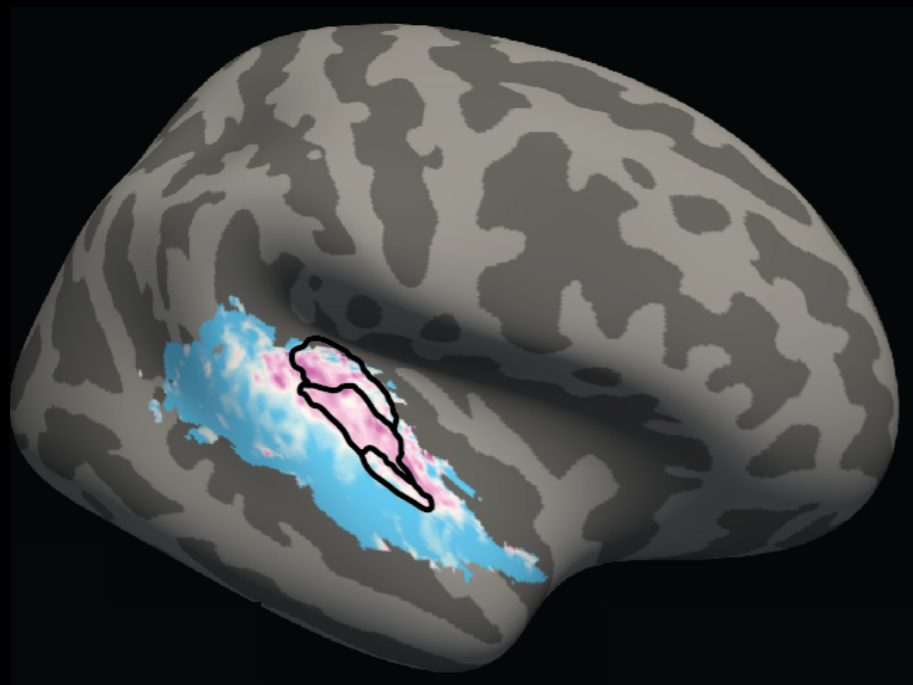**Introduced multi-task networks as neural models**



r²: 0.92

**Performs as well as humans, with similar pattern of errors**

Introduced multi-task networks
as neural models



Performs as well as humans,
with similar pattern of errors



Reveals hierarchical organization
in human auditory cortex

75

# Thanks.



**Josh McDermott**

**Dan Yamins**

**Erica Shook**