

Mind the gap: comparing multiple models of scene representation in brain and behavior

Iris Groen, PhD

New York University

Deep Neural Networks: A New Framework for Modeling Biological Vision and Brain Information Processing

Nikolaus Kriegeskorte

Medical Research Council Cognition and Brain Sciences Unit, University of Cambridge,
Cambridge CB2 7EF, United Kingdom; email: nikolaus.kriegeskorte@mrc-cbu.cam.ac.uk

Neural Networks and Neuroscience-Inspired Computer Vision

David Daniel Cox^{1,2,3,*} and Thomas Dean^{4,5}

Visual Object Recognition: Do We (Finally) Know More Now Than We Did?

Isabel Gauthier¹ and Michael J. Tarr²

Using goal-driven deep learning models to understand sensory cortex

Daniel L K Yamins^{1,2} & James J DiCarlo^{1,2}

Toward an Integration of Deep Learning and Neuroscience

Adam H. Marblestone^{1*}, Greg Wayne² and Konrad P. Kording³

The 5th CiNet Conference

Computation and representation in brains and machines



February 20-22, 2019

Center for Information and Neural Networks, Osaka, Japan

Meeting Chair: Shoji Nishimoto (CiNet/NICT)

Ce-chair: Shigeru Kitazawa (CiNet/Osaka University)

Takafumi Suzuki (CiNet/NICT)

Meeting Director: Takahisa Taguchi (CiNet/NICT)

Organizer:

NICT, Center for Information and Neural Networks

Sponsors:

-Grant-in-Aid for Scientific Research on Innovative Areas, MEXT, Japan

-“Chronogenesis: How the Mind Generates Time”

-NEC Corporation

-NTT Data Institute of Management Consulting, Inc.

Financial support: Ichimura Foundation for New Technology

Program and registration details

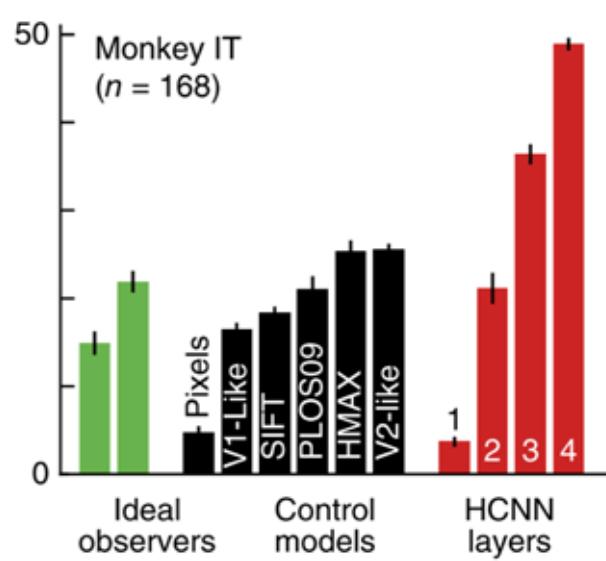
https://cinen.jp/english/event/5th_cinetconf/

Speakers:
Shun-ichi Amari
Matthew Botvinick
David Cox
Dileep George
Marcel van Gerven
Iris Groen
Michael Hainke
Uli Hasson
Aspo Hyvärinen
Yukiyasu Kamitani
Shigeru Kitazawa
Nikolaus Kriegeskorte
Jun Morimoto
Tomoya Nakai
Satomi Nishida
Shoji Nishimoto
Ana Luisa Pinho
Odeila Schwartz
Taro Toyozumi
Kai Wang
Daniel Yamins
Takafumi Yanagisawa

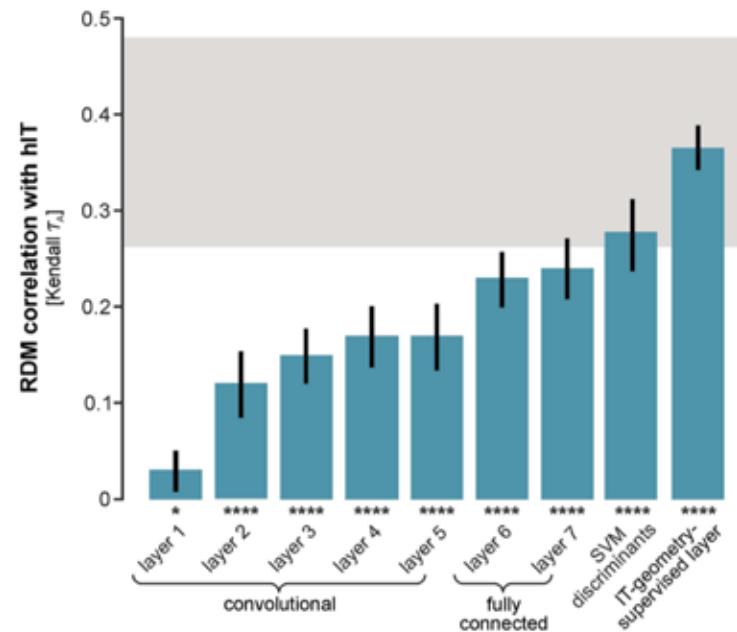
Deep nets and object recognition



Krizhevsky,
Sutskever
& Hinton, 2012



Yamins et al, 2014



Khaligh-Razavi
& Kriegeskorte, 2014

beach hut
sand
beach chair

“beach”

Does deep learning explain scene perception?



Outline

- Scene vs. object perception
- fMRI study 1: objects-in-context
- fMRI study 2: comparing multiple models
- How to move forward?

Scene vs. object perception

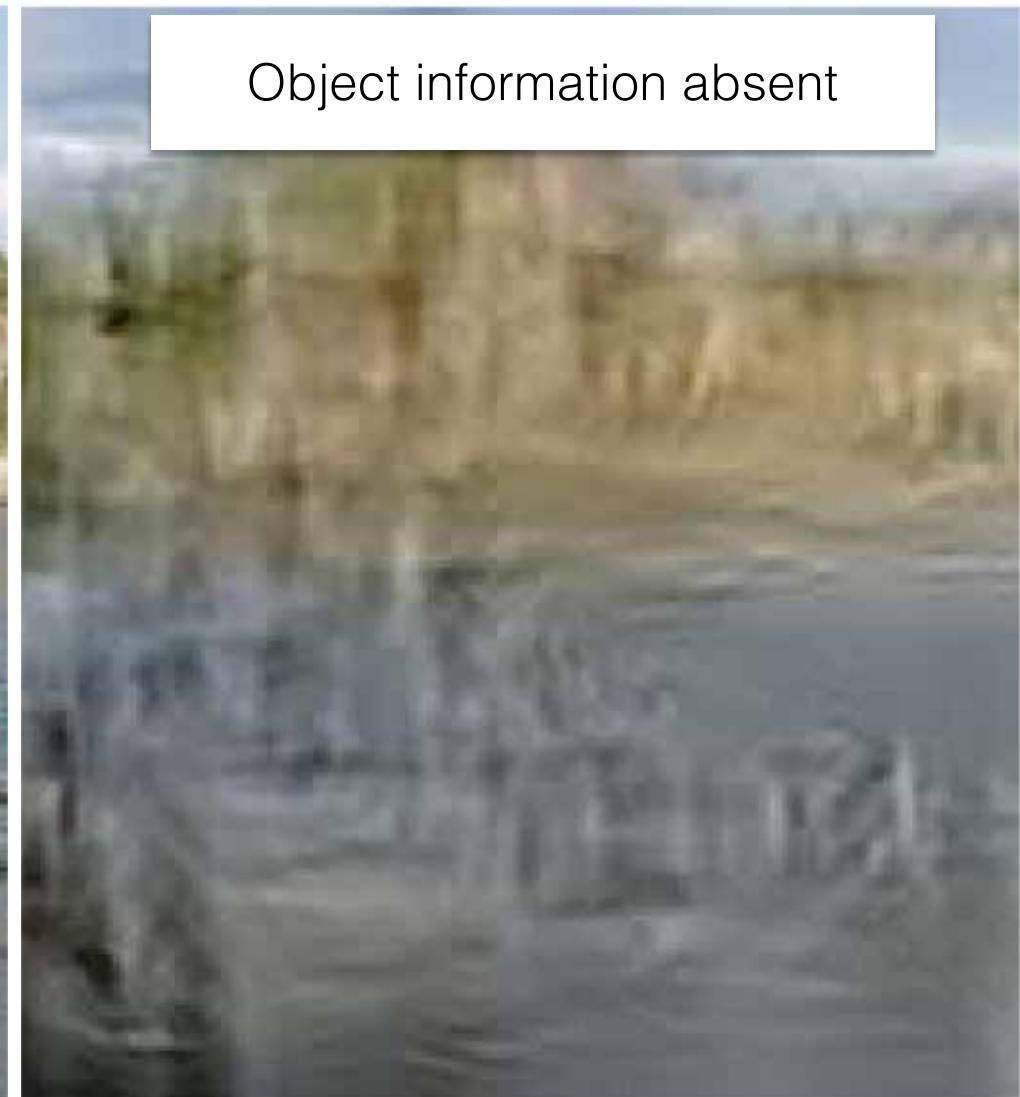
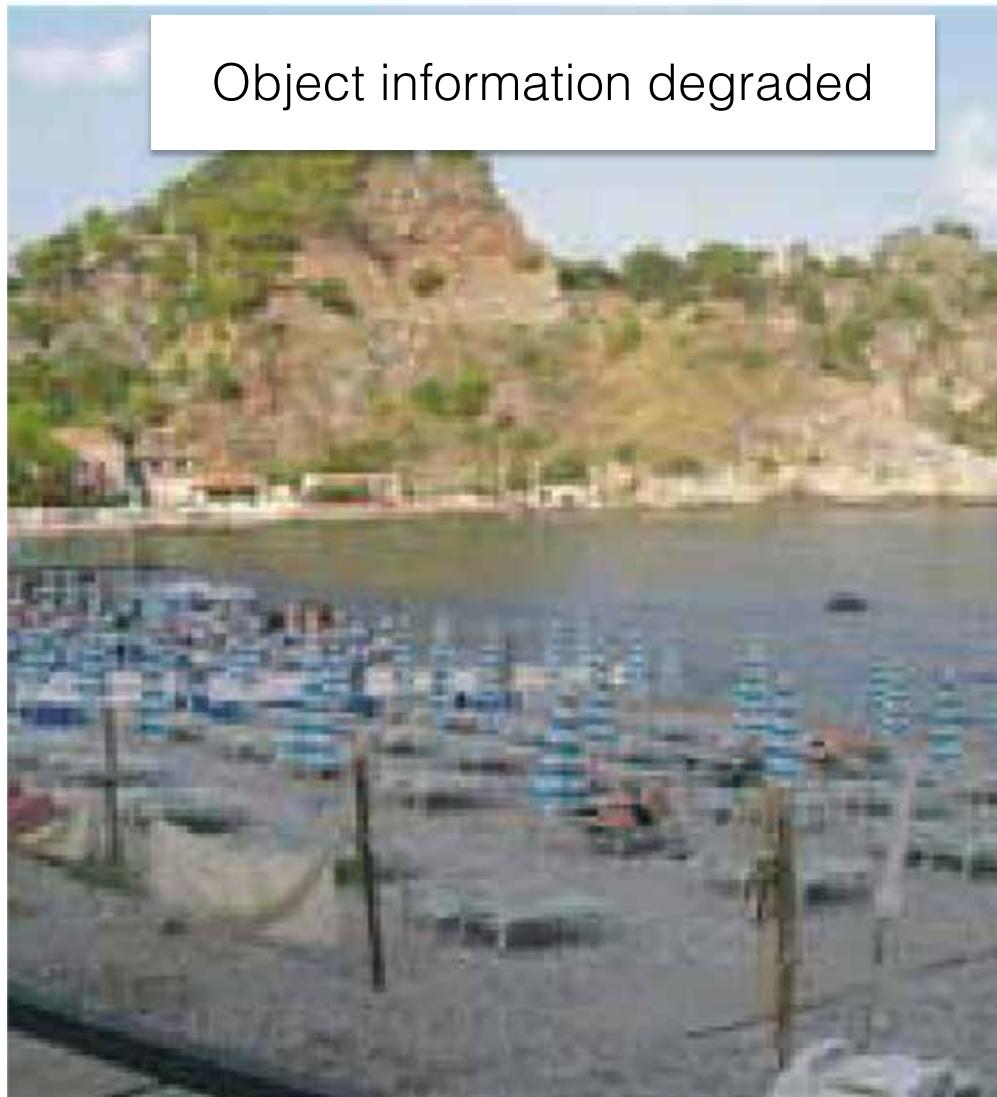
A scene is a semantically coherent (and often namable) view of a real-world environment comprising **background elements** and multiple discrete objects arranged in a **spatially licensed manner**.

Henderson and Hollingworth (1999)

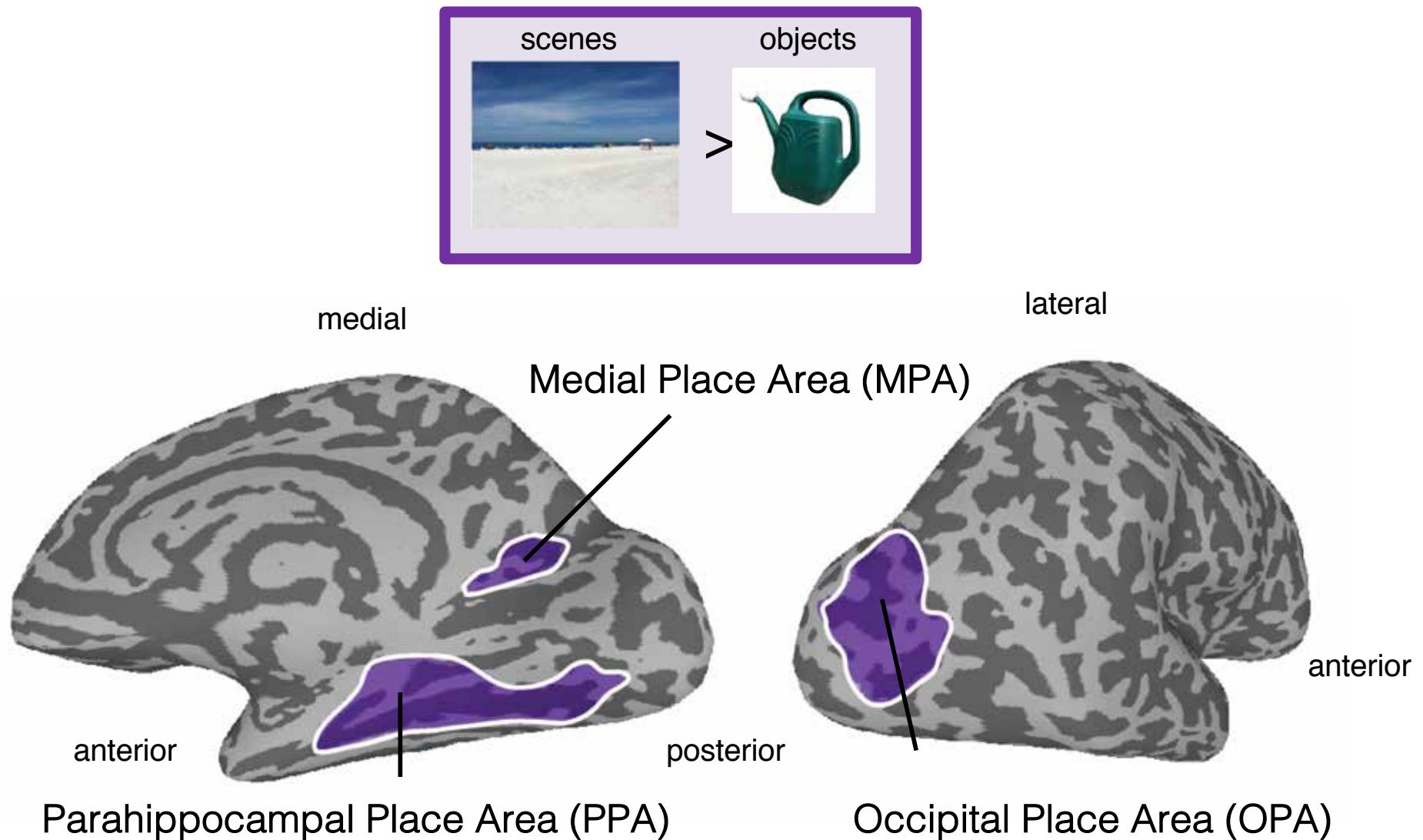
Scene perception can be usefully contrasted to object perception: whereas objects are **spatially compact entities** that one acts upon, scenes are **spatially distributed entities** that one acts within.

Epstein (2005)

Scene recognition without objects

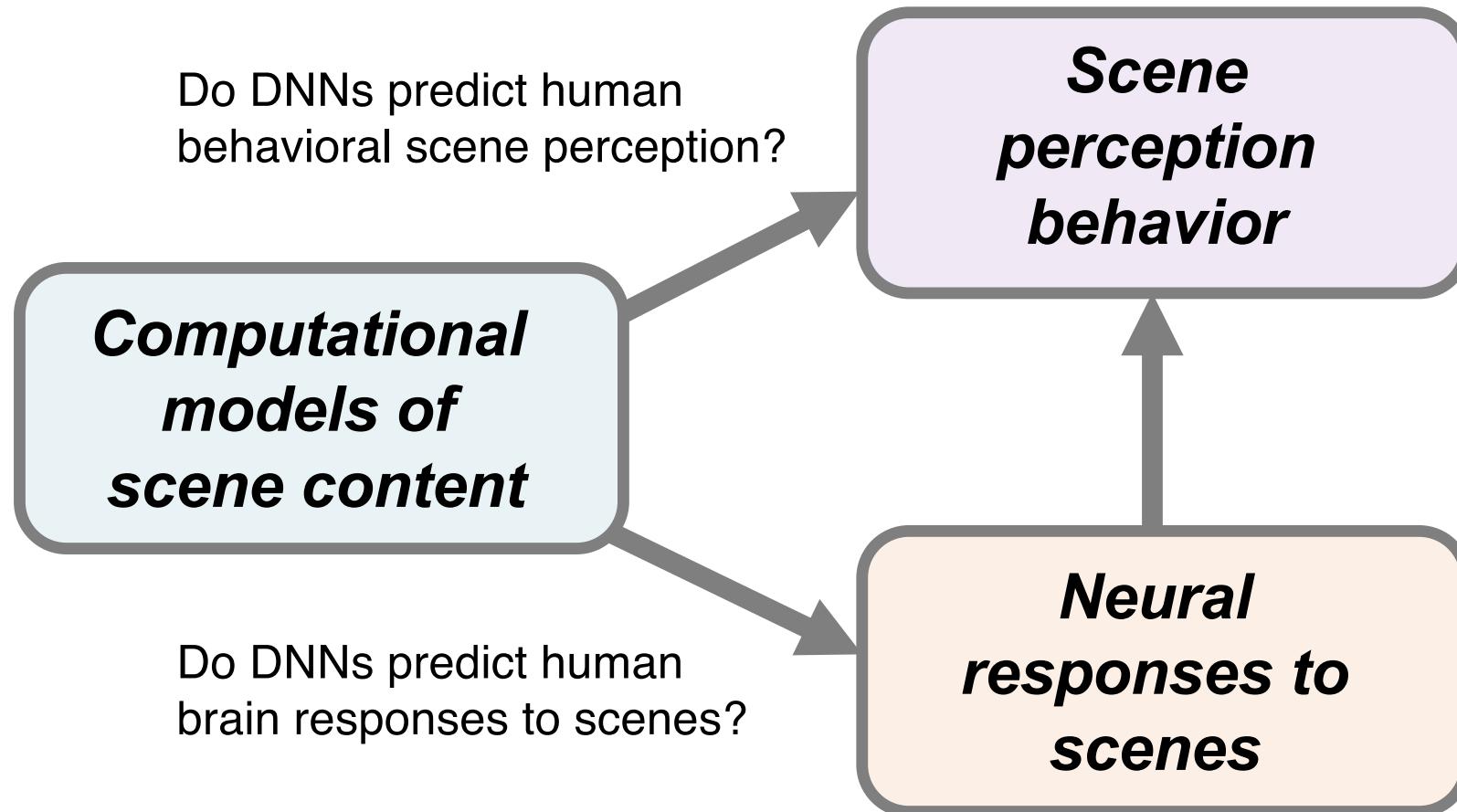


Scene-selective brain regions



See e.g. Epstein (2014) for a review

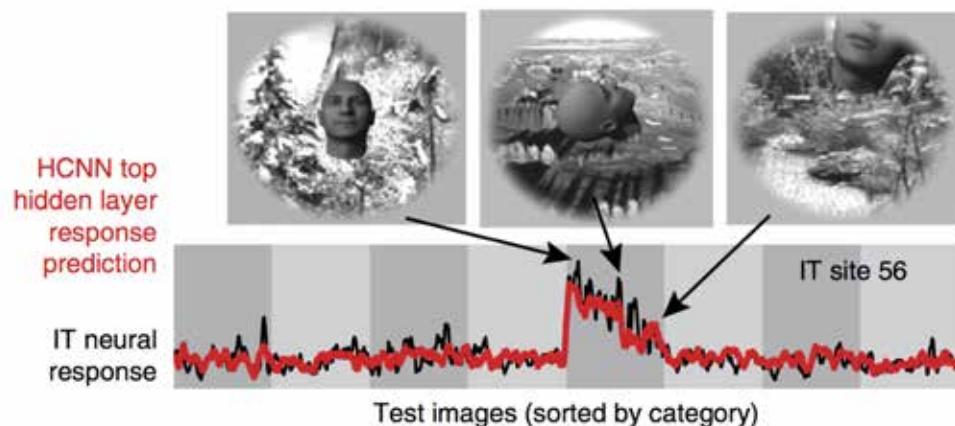
Understanding scene perception



Outline

- Scene vs. object perception
- fMRI study 1: objects-in-context
- fMRI study 2: comparing multiple models
- How to move forward?

DNNs and object recognition

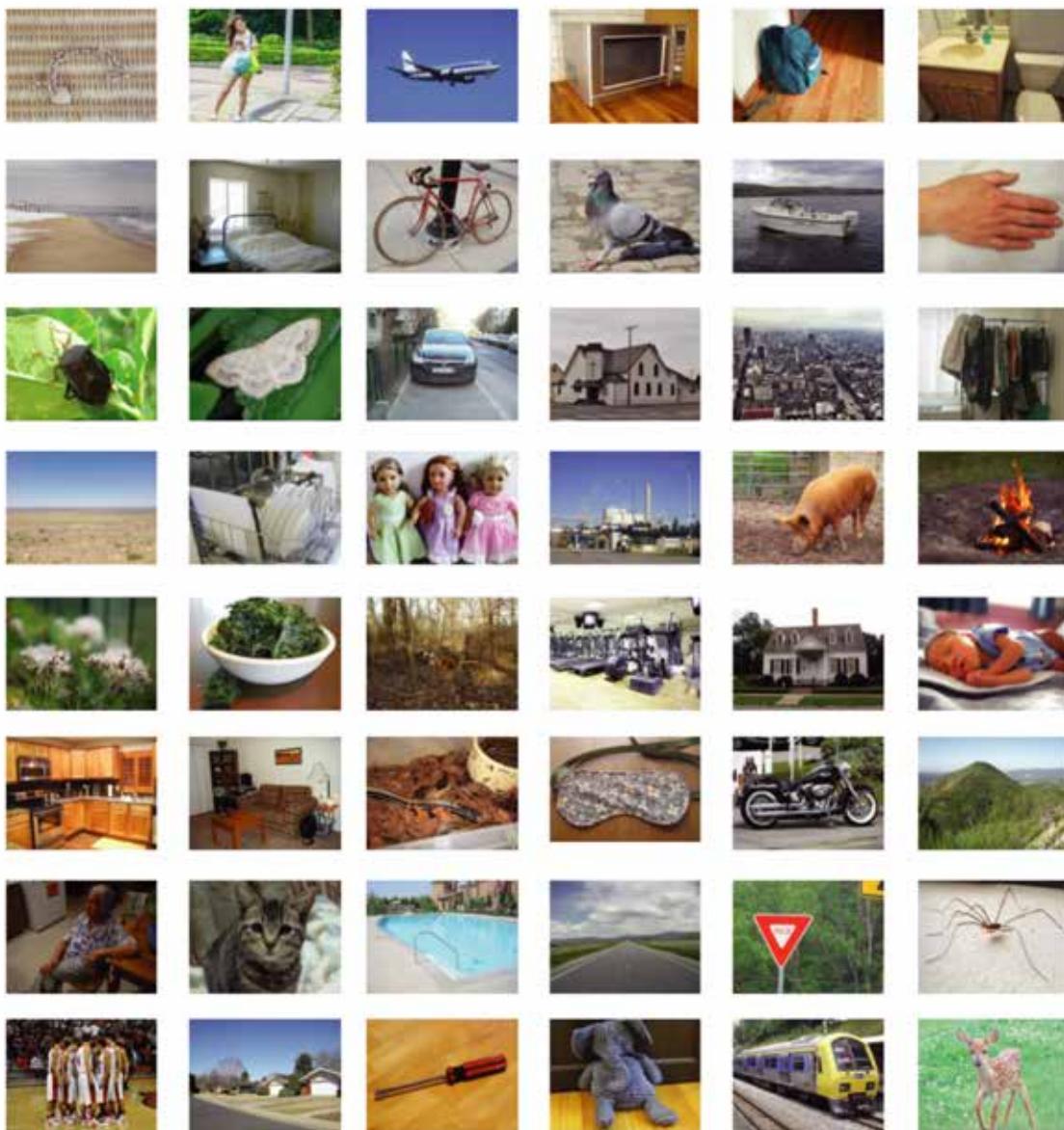


Yamins et al, 2014



*Khaligh-Razavi
& Kriegeskorte, 2014*

Study 1: objects-in-context + scenes

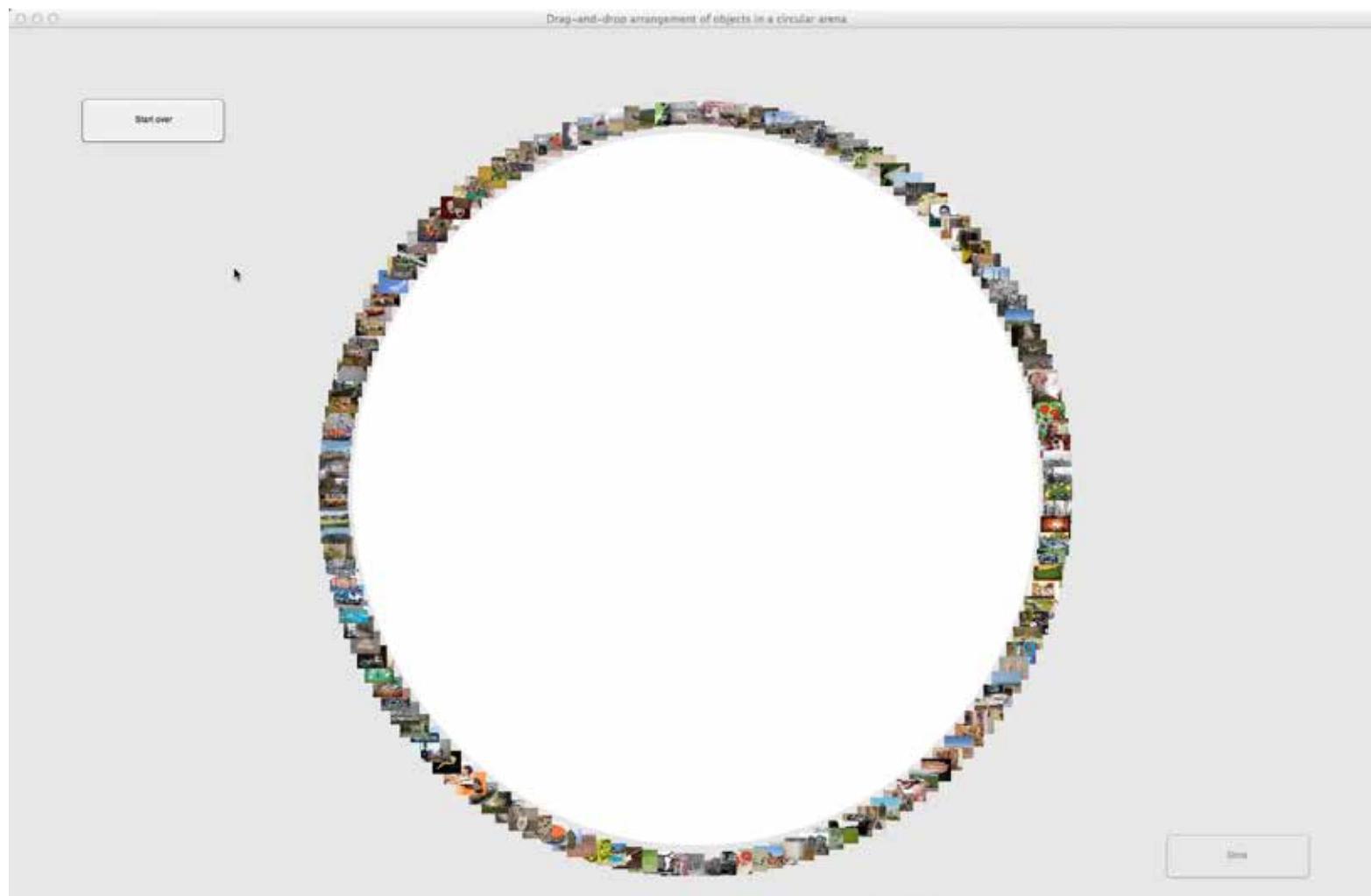


**48 naturalistic image categories
(3 exemplars per category, 2 sets)**

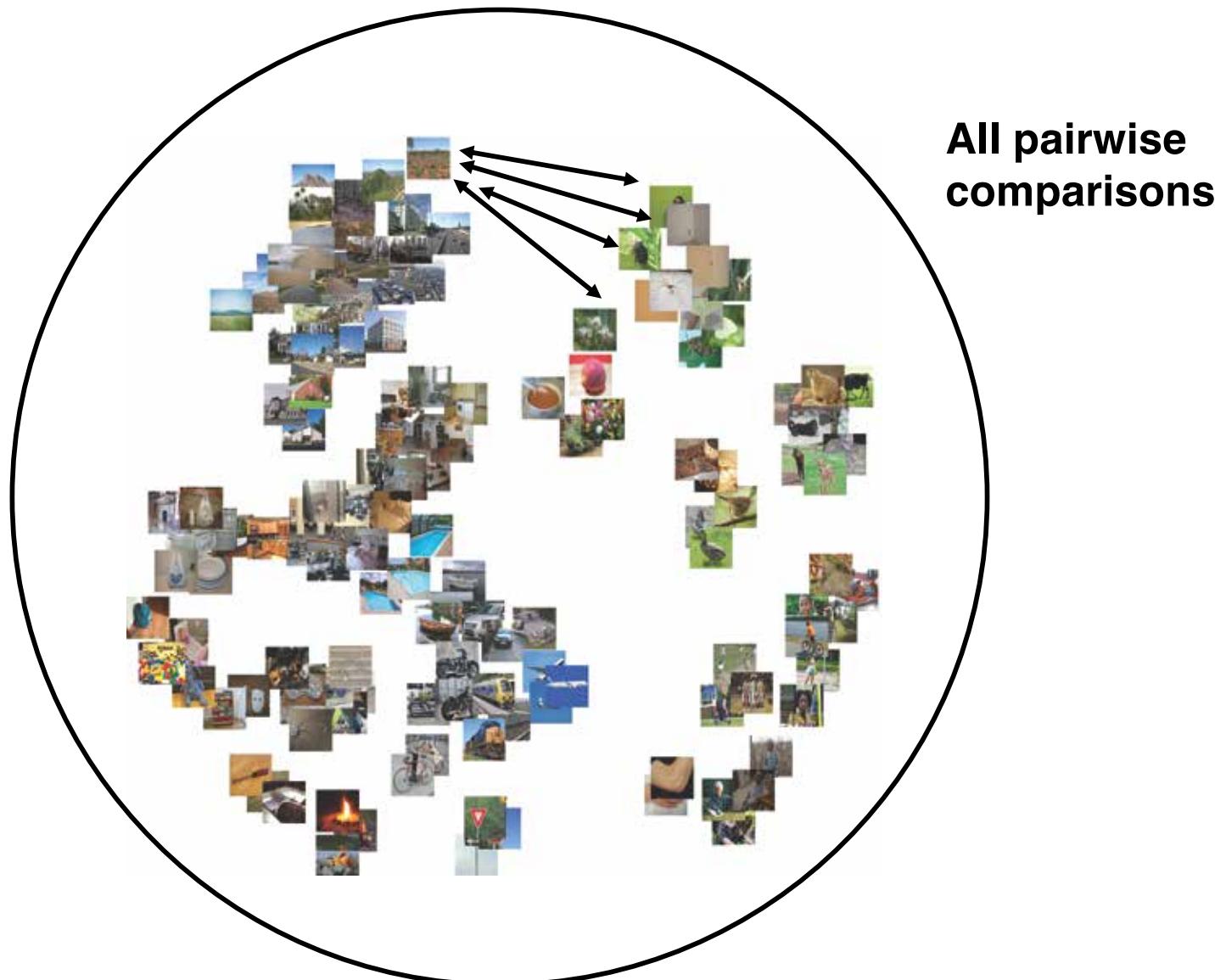
accessories	flowers
adults	food
airplanes	forests
appliances	gyms
bags	houses
bathrooms	kids
beaches	kitchens
beds	living rooms
bikes	lizards/snakes
birds	masks
boats	motorcycles
body parts	mountains
bugs	older adults
butterflies	pets
cars	pools
churches	roads
cityscapes	signs
clothes	spiders
deserts	sports
dishes	suburbs
dolls	tools
factories	toys
farm animals	trains
fire	wild animals

Free arrangement behavioral task

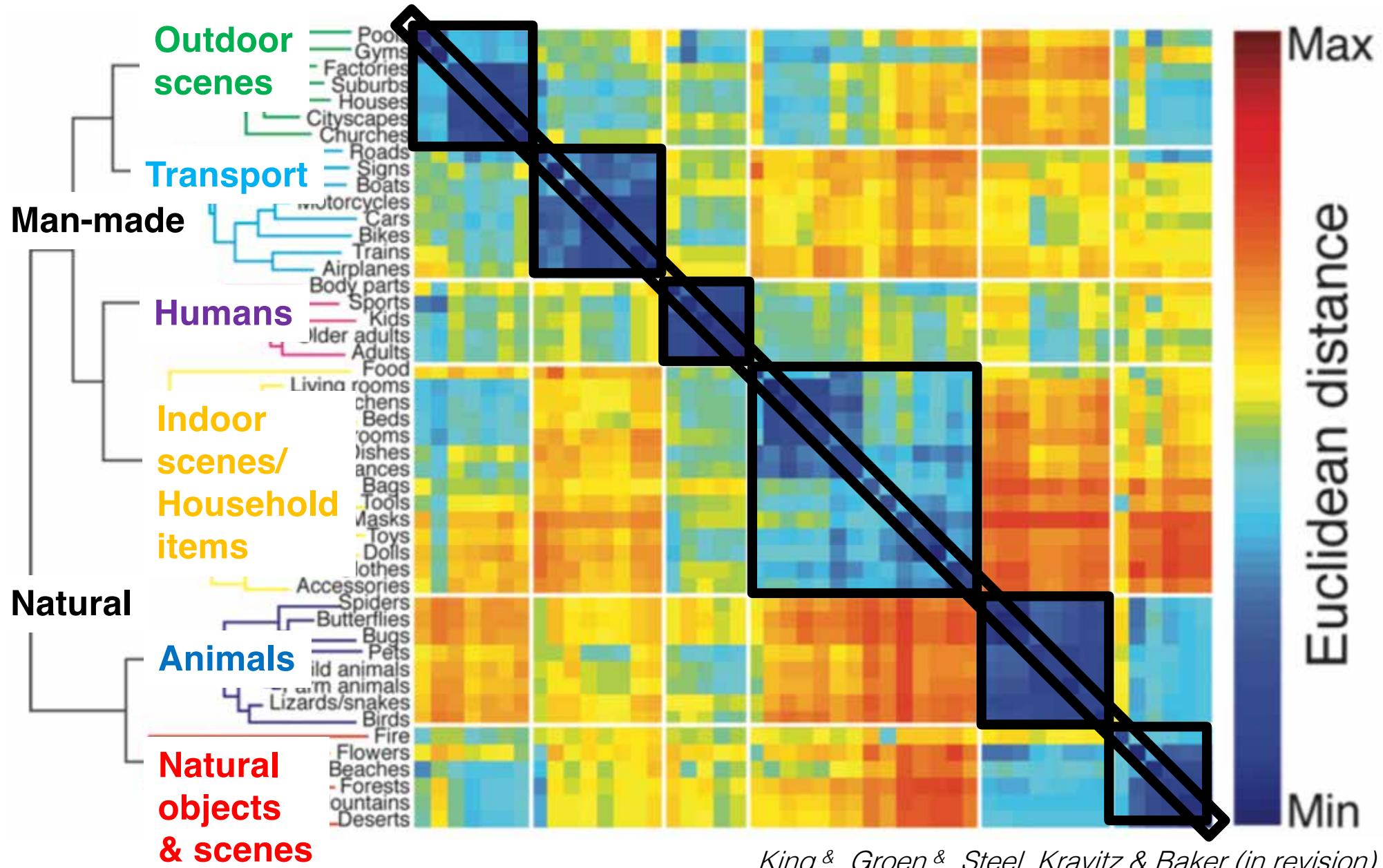
“Group these images according to their similarity”



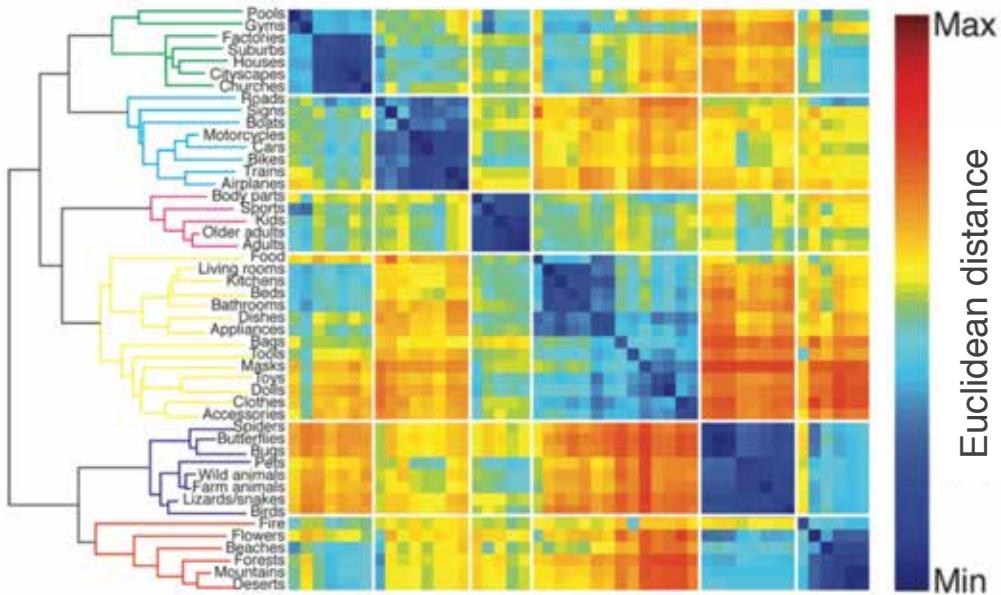
Free arrangement behavioral task



Behavioral dissimilarity



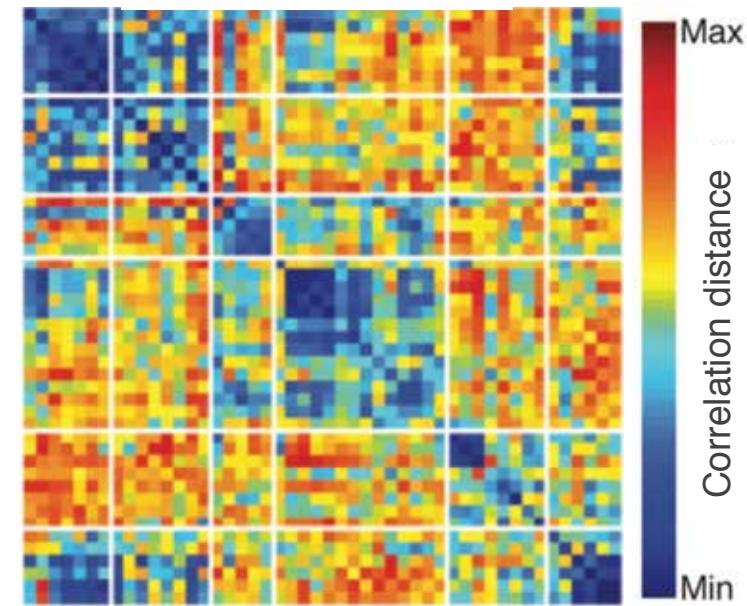
Behavior



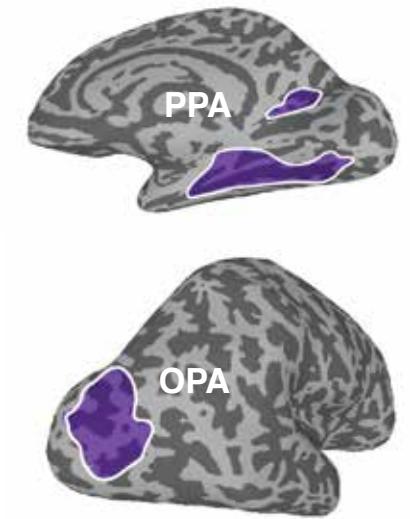
Ventral temporal cortex?



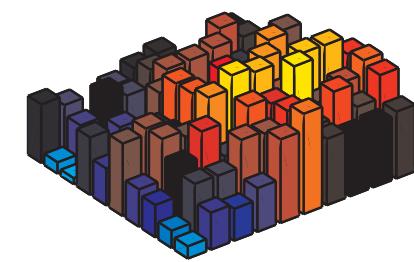
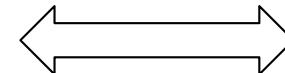
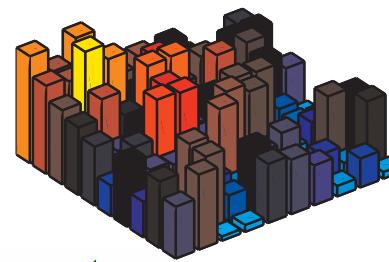
AlexNet layer 8



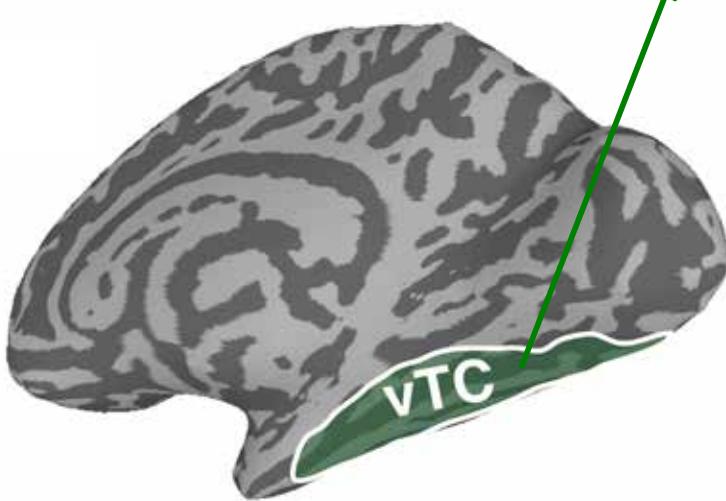
Scene-selective cortex?



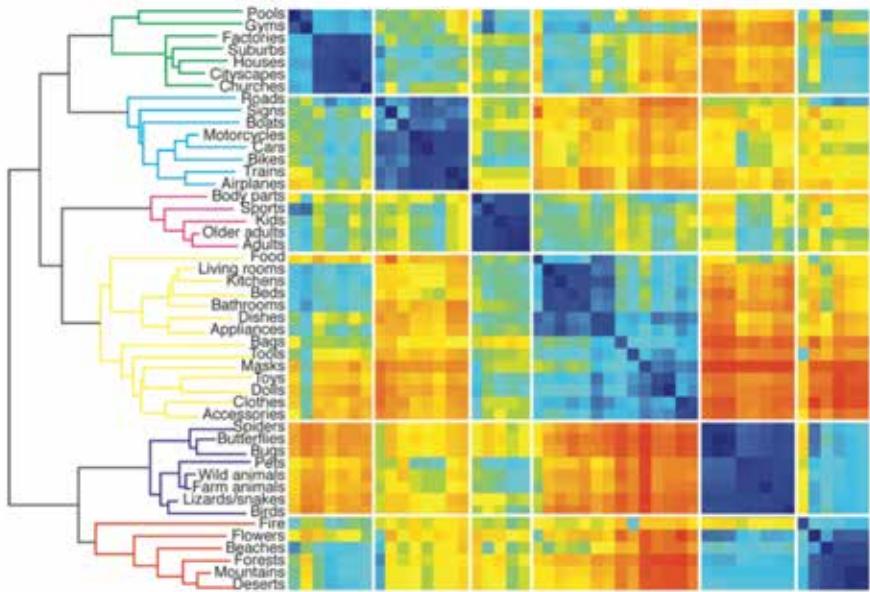
Multi-voxel pattern analysis



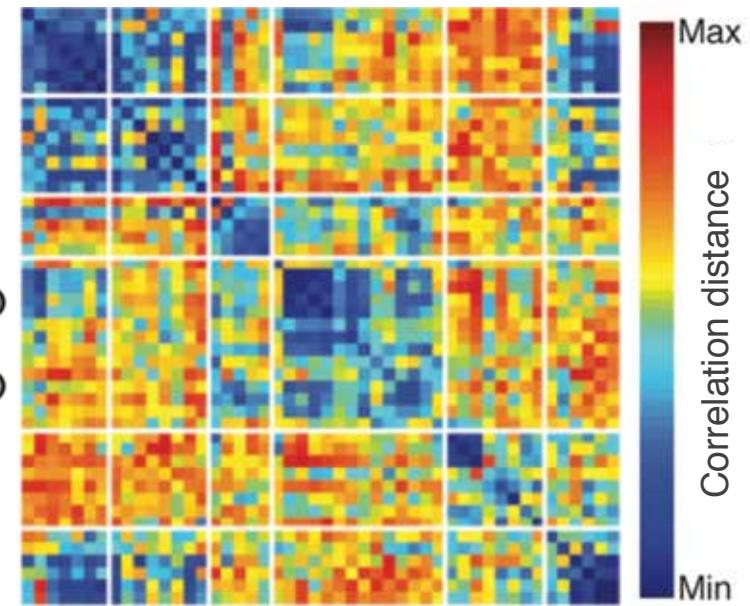
1-correlation



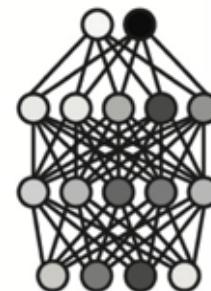
Behavior



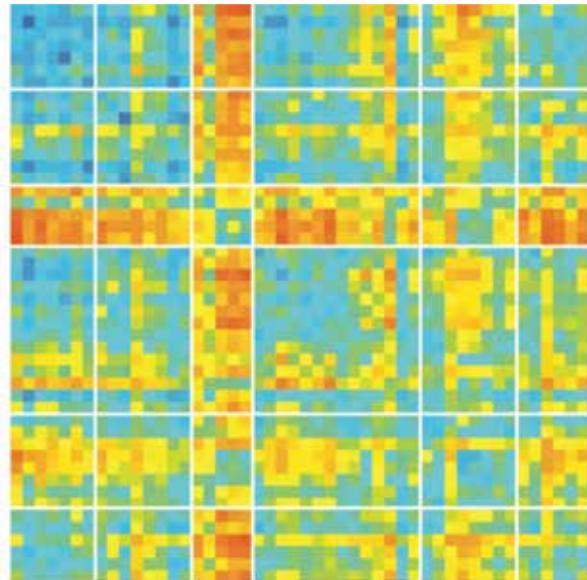
AlexNet layer 8



DNN

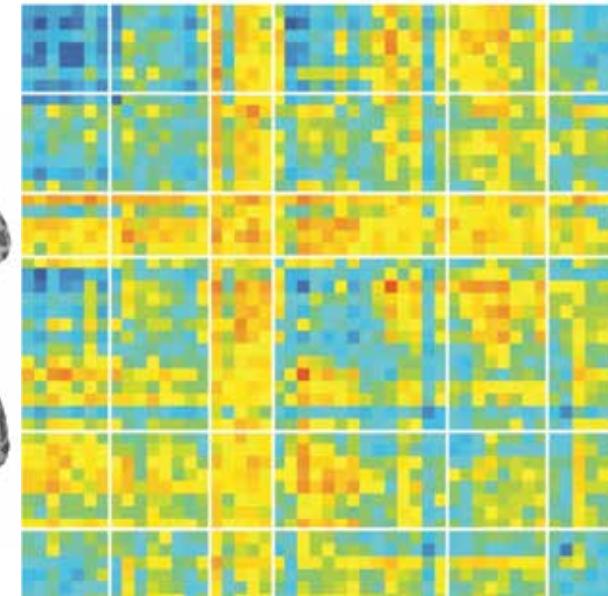
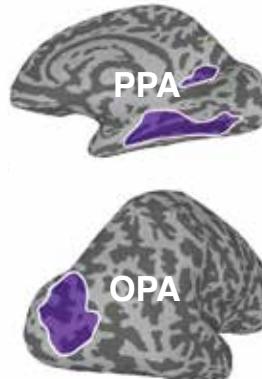


Ventral temporal cortex

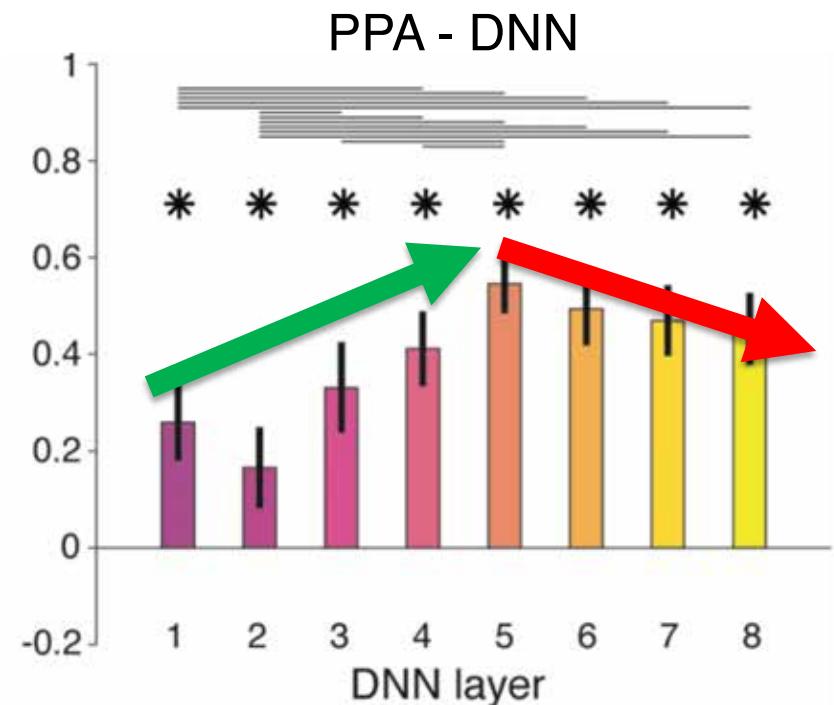
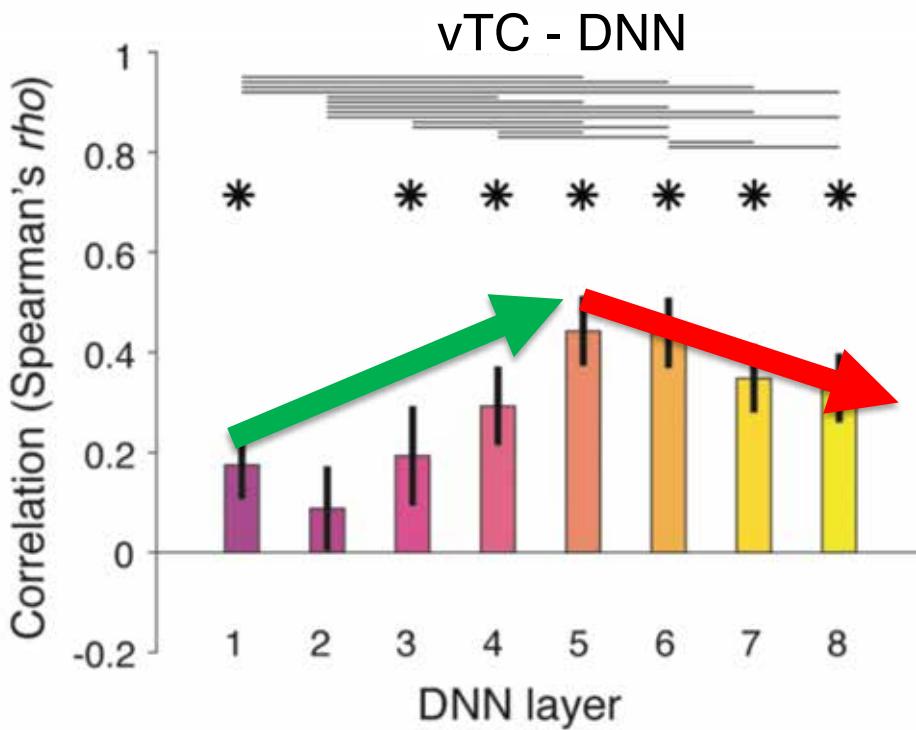
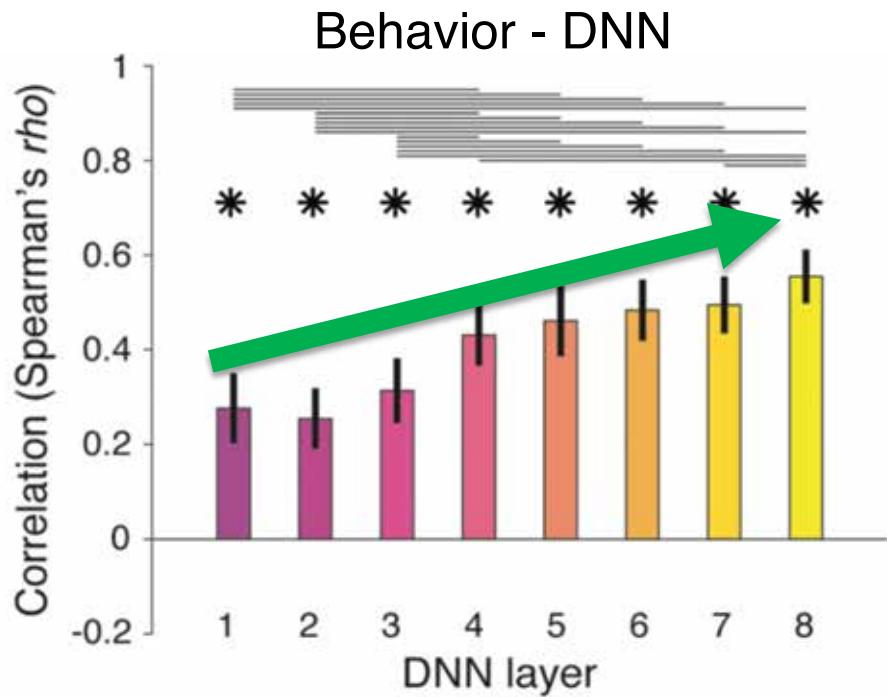


Max
Correlation distance
Min

Scene-selective cortex



Max
Correlation distance
Min



Understanding scene perception

Do DNNs predict human behavioral scene perception?

Yes - layer 8 best

Computational models of scene content

Do DNNs predict human brain responses to scenes?

Yes - layer 5 best

Scene perception behavior

Mismatch!

Neural responses to scenes



Why the mismatch?

- Nature of stimuli (mix of objects-in-context and scenes)?
- Multi-arrangement task?
- Alternative models for behavior?

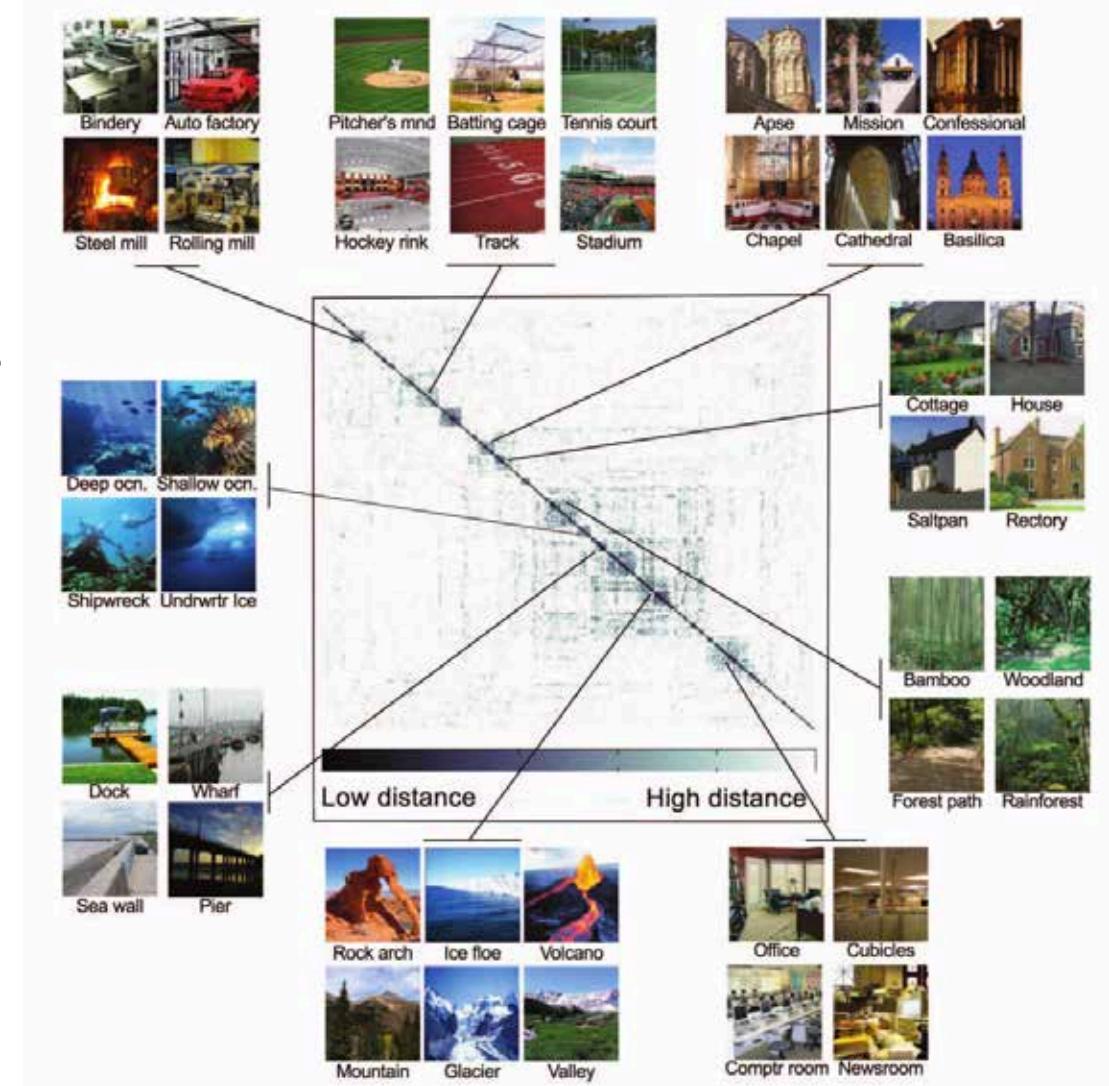
Large-scale scene perception

- Mechanical Turk
- SUN database
- > 300 scene categories

SAME or DIFFERENT CATEGORY?

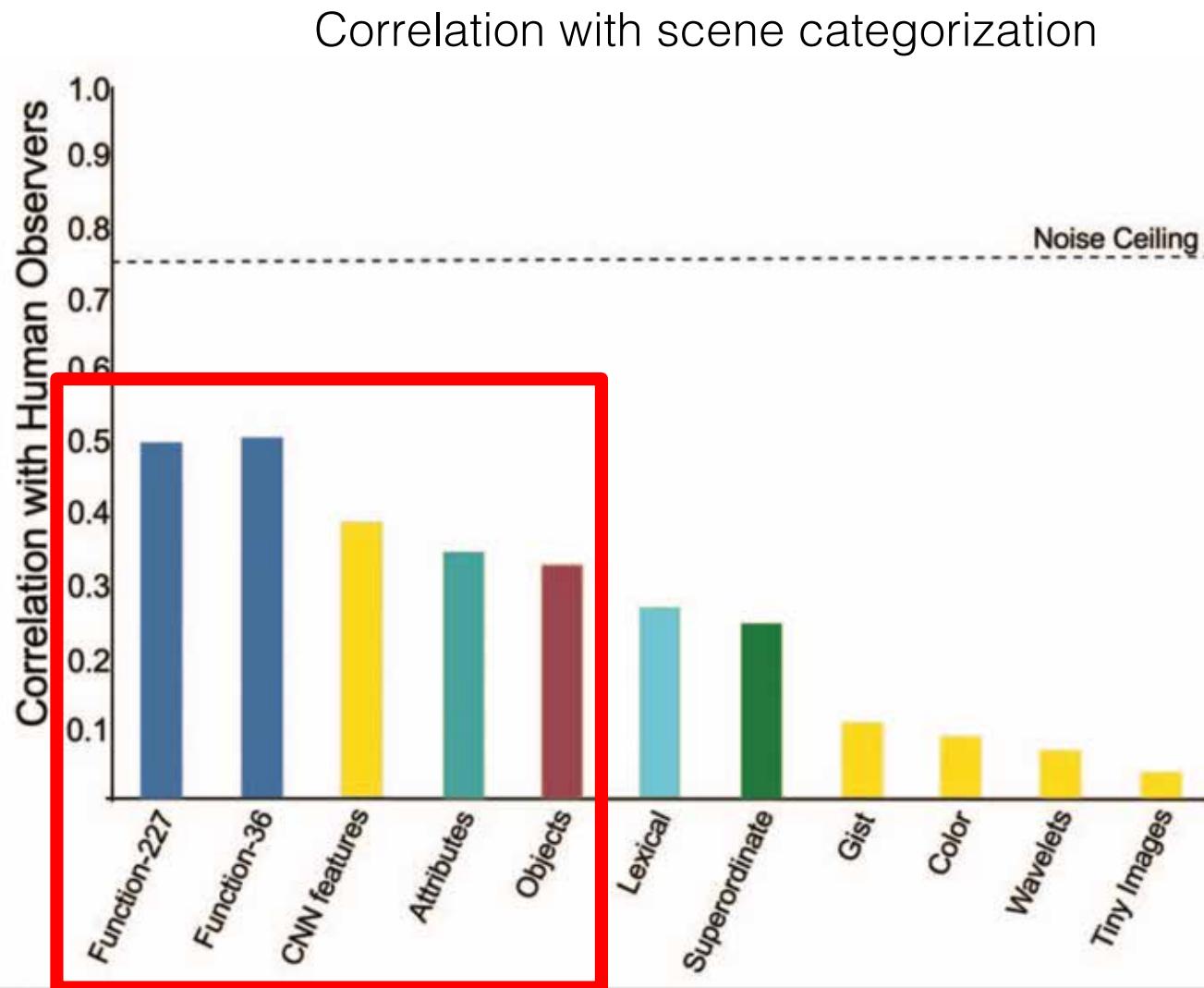


> 2000 participants
> 5 million trials



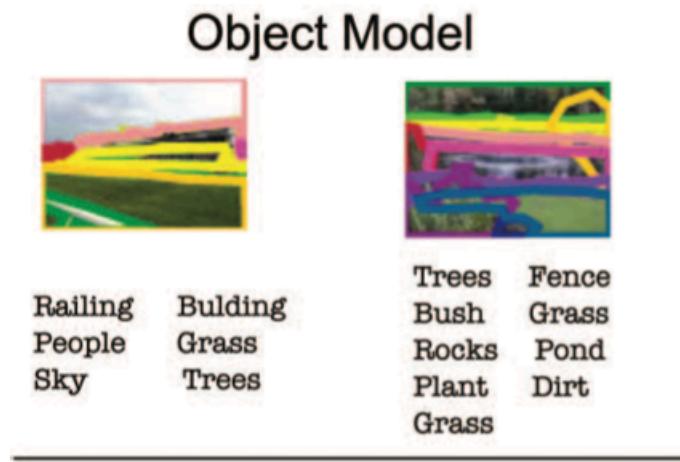
Greene et al., 2016 (JEP:General)

Comparing behavior to multiple models



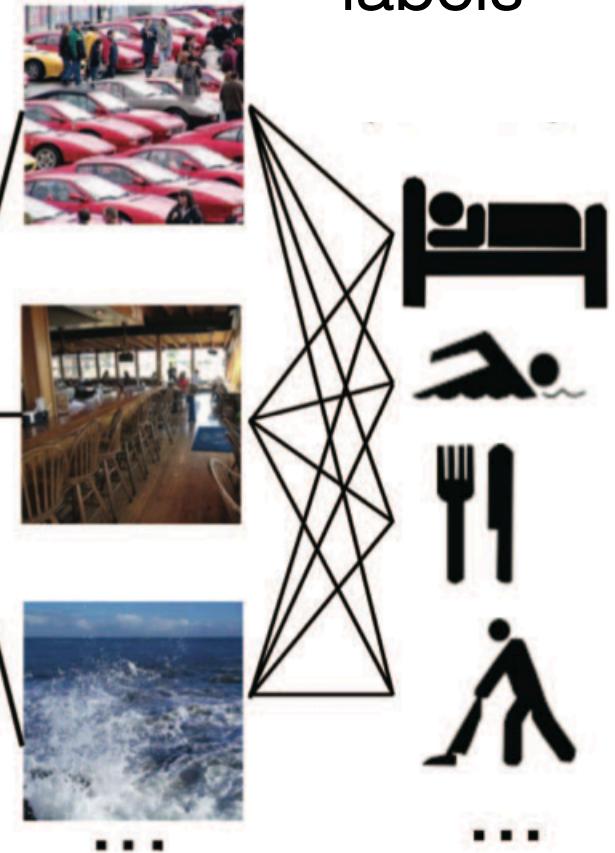
Greene et al., 2016 (JEP:General)

Three best models of scene perception



Function Model

Scenes **Action labels**



Perceptual Model



Convolutional Neural Network
(AlexNet; 7th layer)

Greene et al., 2016 (JEP:General)

Study 2

- **Goal:** determine how well the top three models predict fMRI responses and behavioral multi-arrangement task

!! Inherent correlations make it difficult to disambiguate models, even within large sets of naturalistic images

(Lescroart, Stansbury & Gallant, 2015; Malcolm, Groen & Baker, 2016)

- Minimize model correlations through stimulus selection:
Iterative sampling from SUN database until we obtain maximally different predictions
- Use variance partitioning to identify unique contribution of each model

Stimuli by function

access_road



airplane_cabin



apse



badminton_court



bamboo_forest



bar



batting_cage



bindery



bus_depot



butte



control_tower



dolmen



escalator



hedgerow



lido_deck



naval_base



pilothouse



playroom



pump_room



putting_green



skyscraper



stilt_house



tea_garden



underwater_pool



volcano



volleyball_court



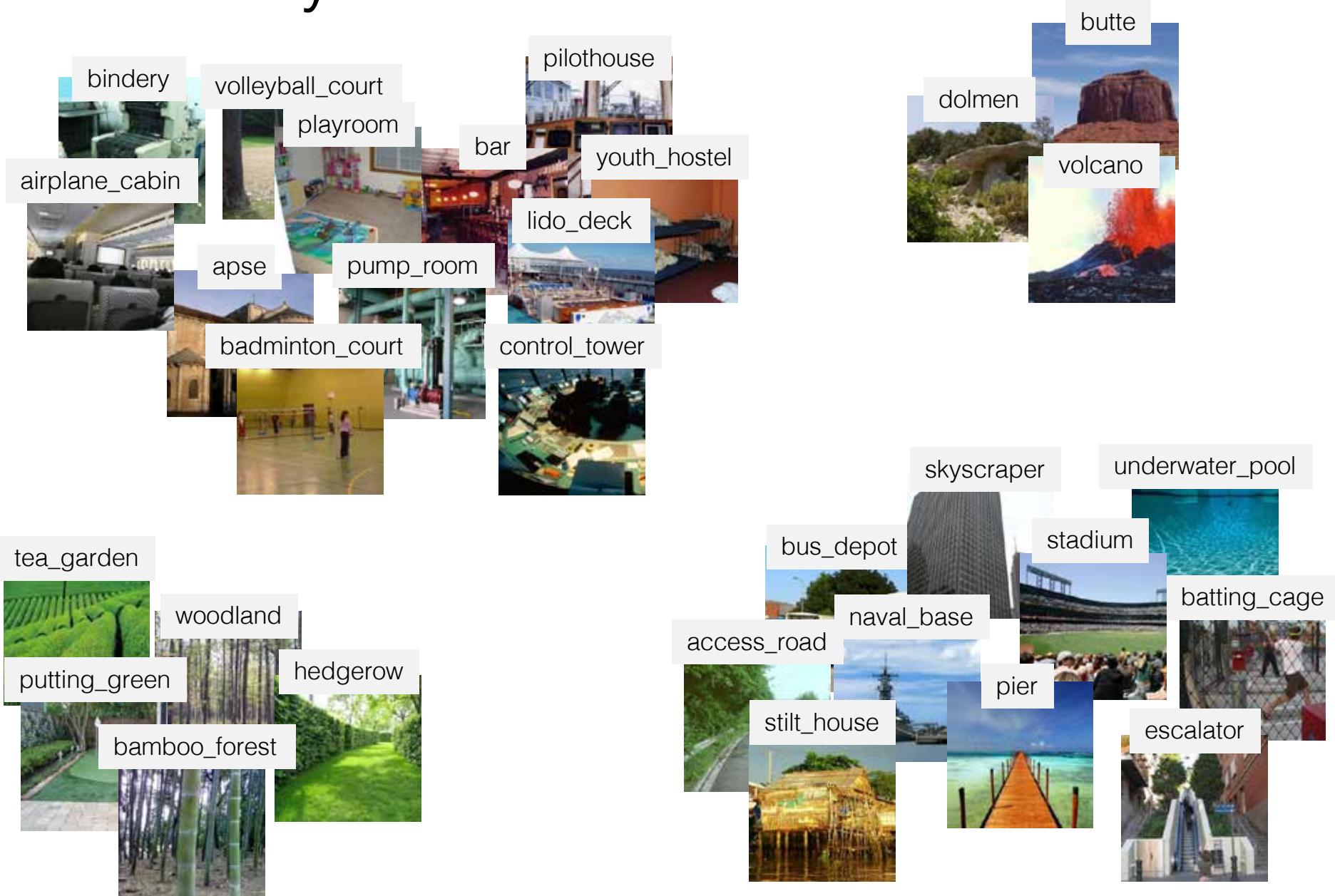
woodland



youth_hostel

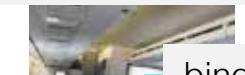


Stimuli by DNN



Stimuli by Objects

airplane_cabin



control_tower



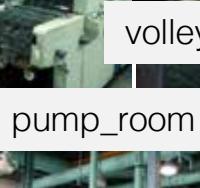
youth_hostel



bar



bindery



volleyball_court



pump_room



underwater_pool



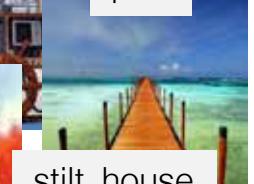
pilothouse



lido_deck



pier



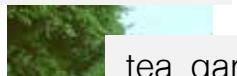
volcano



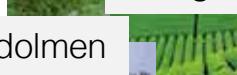
stilt_house



access_road



tea_garden



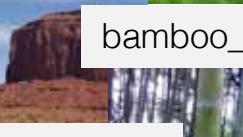
hedgerow



dolmen



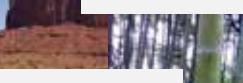
butte



woodland



bamboo_forest



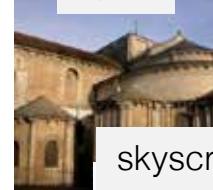
batting_cage



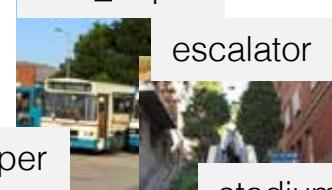
playroom



apse



bus_depot



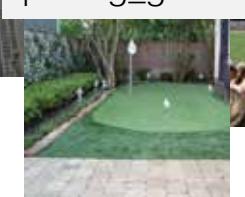
escalator



skyscraper



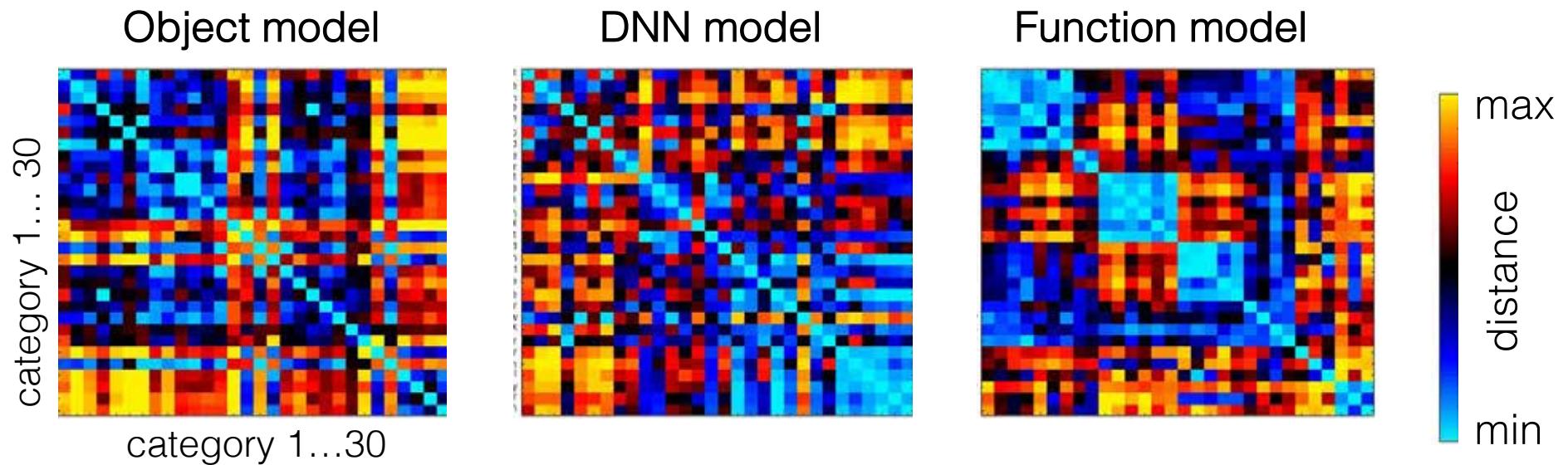
putting_green



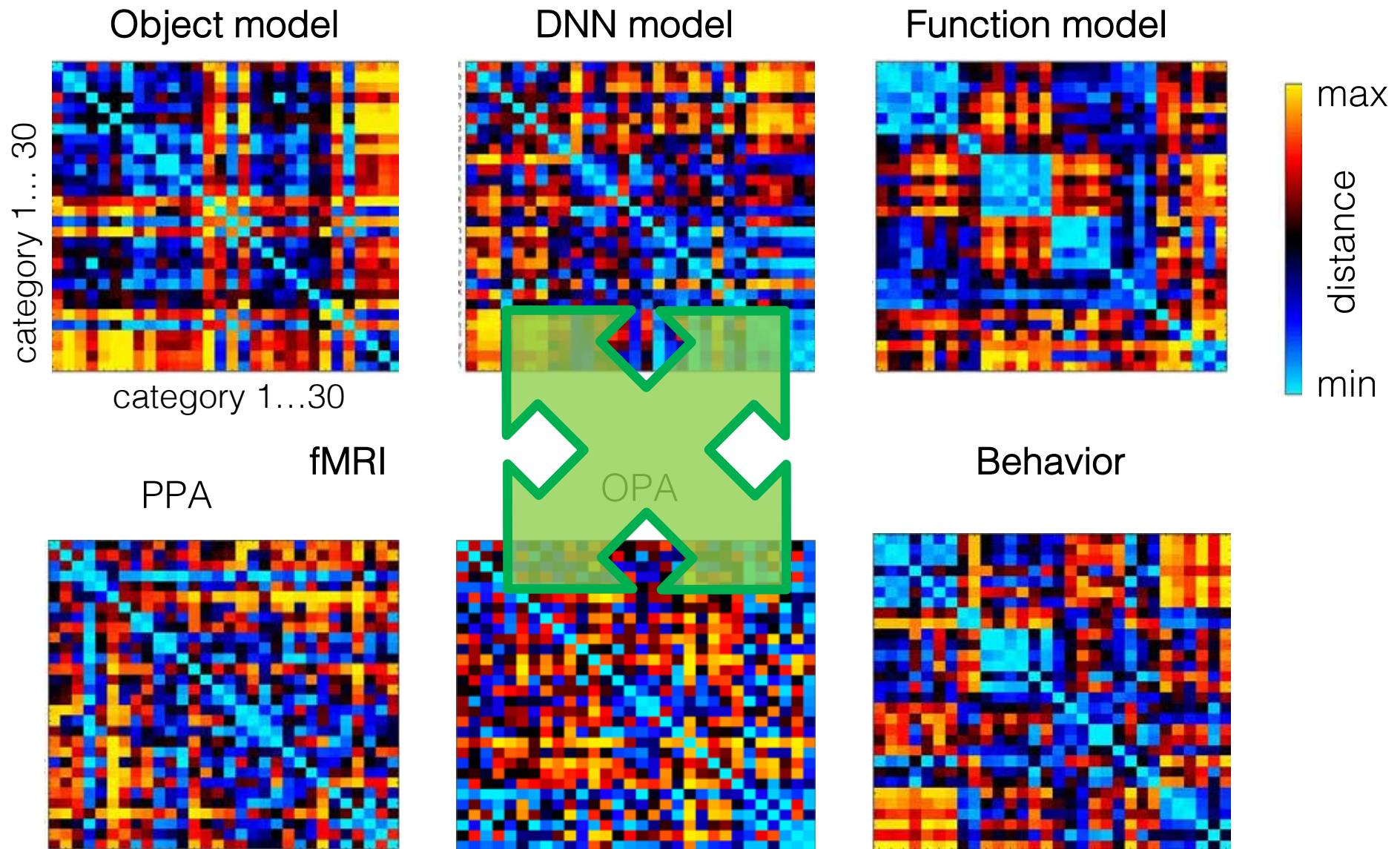
naval_base



Model predictions

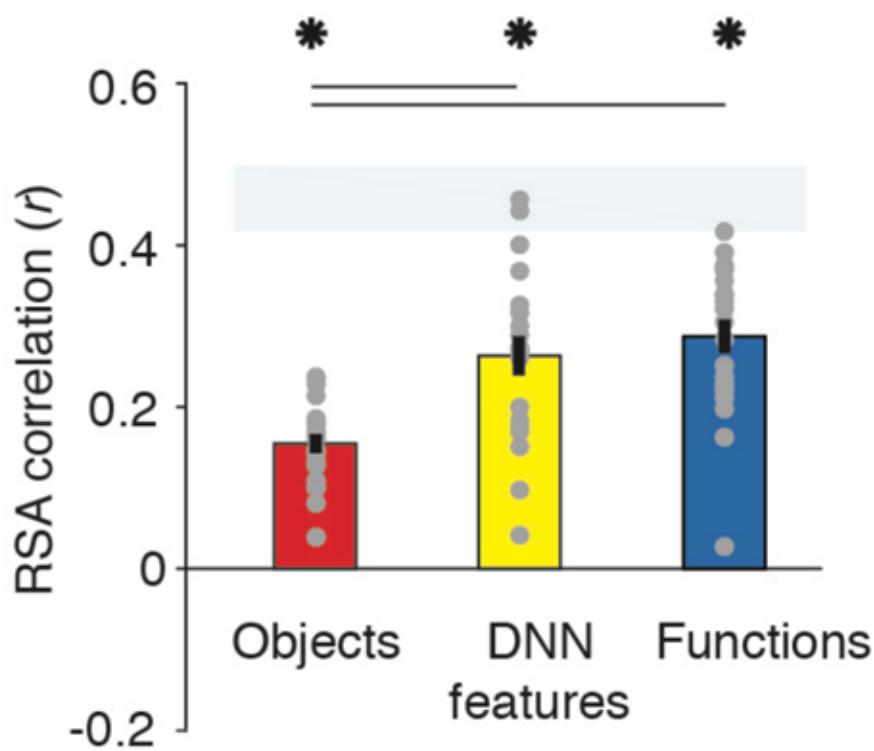


Model predictions

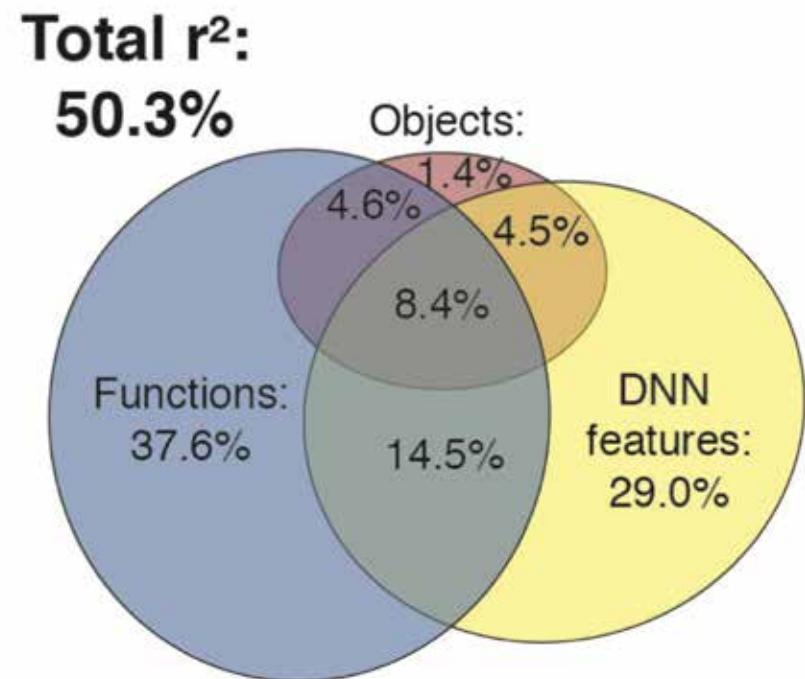


Model correlations with behavior

Correlations

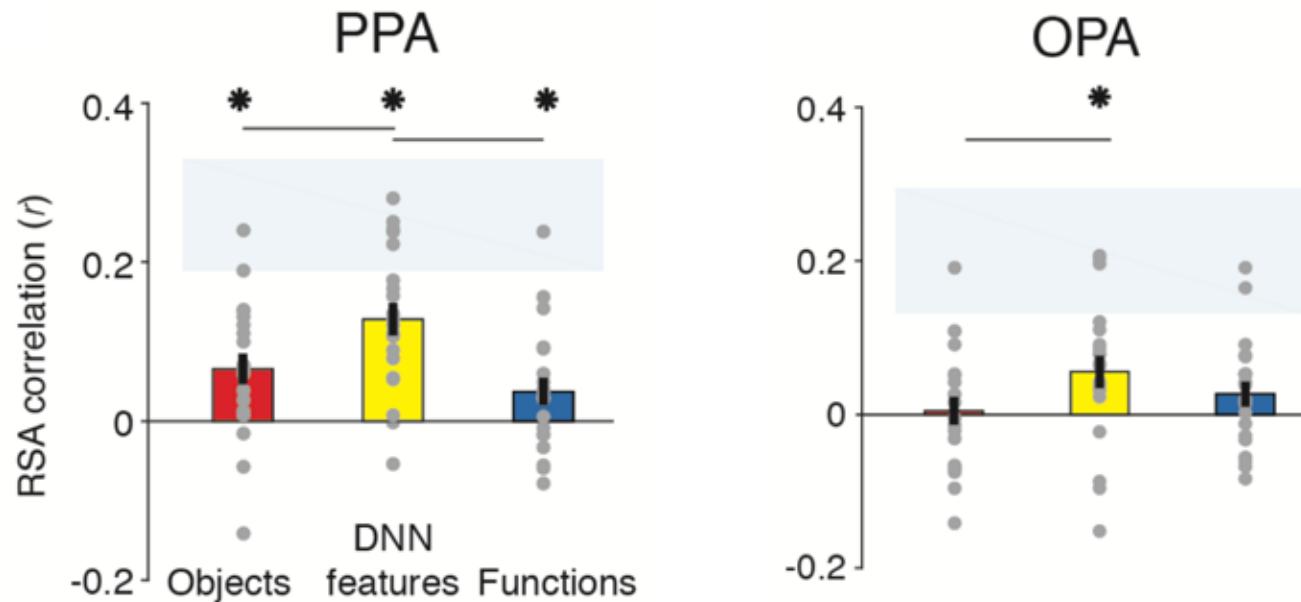


Variance partitioning

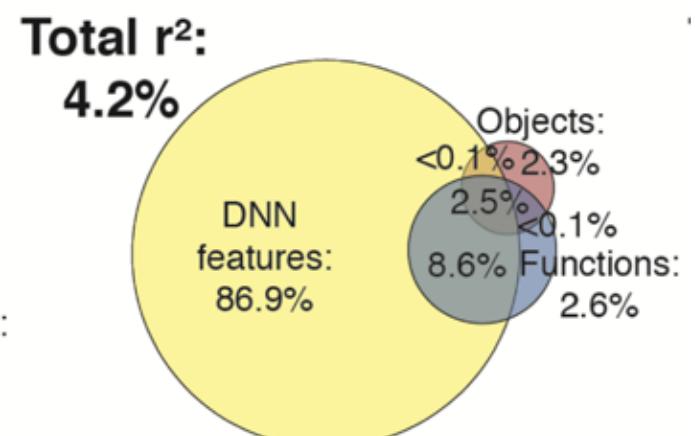
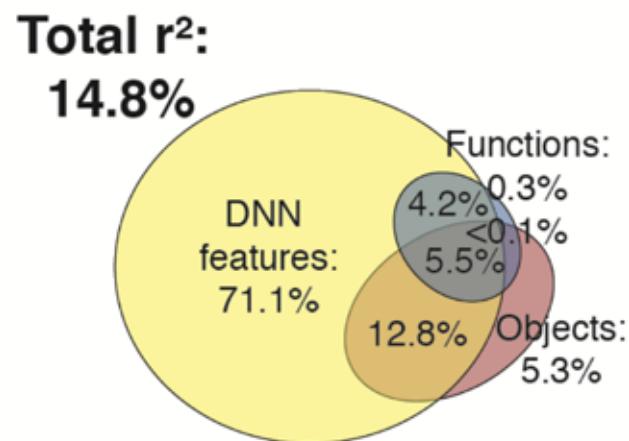


Model correlations with brain

Correlations

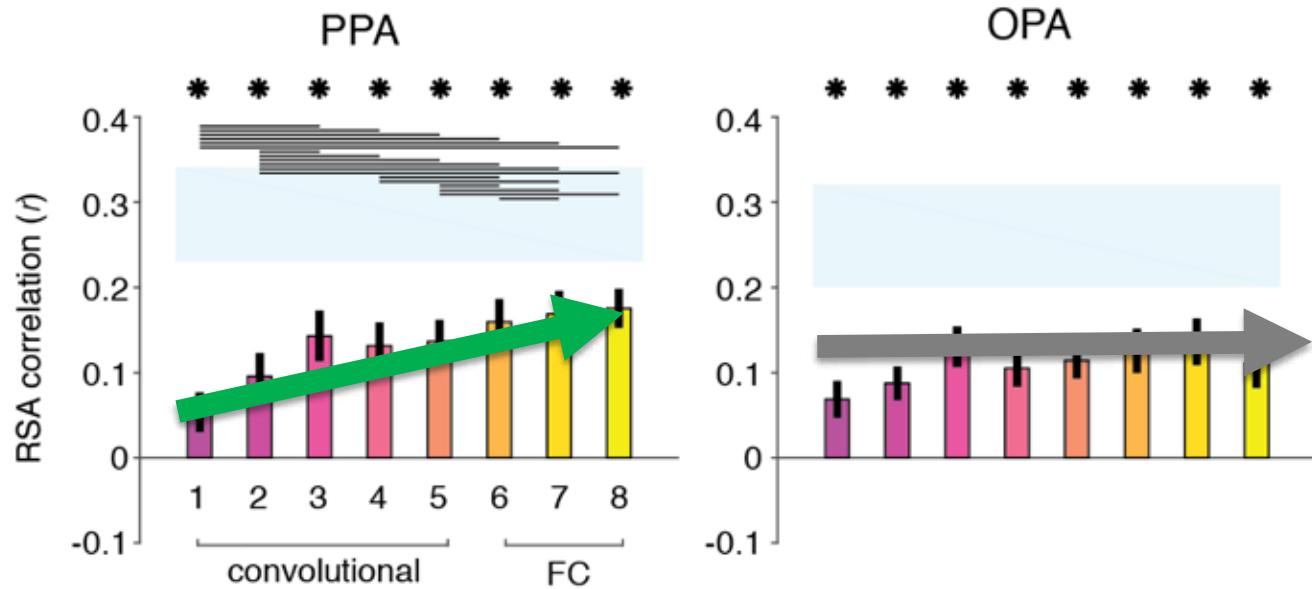


Variance partitioning

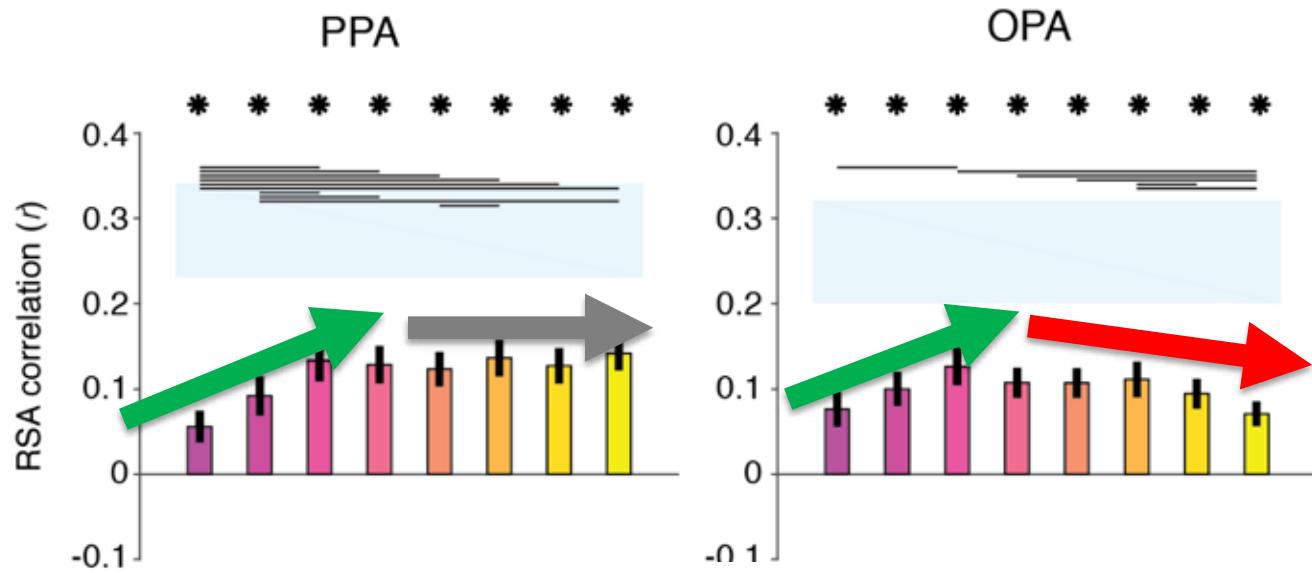


DNN correlations in scene cortex

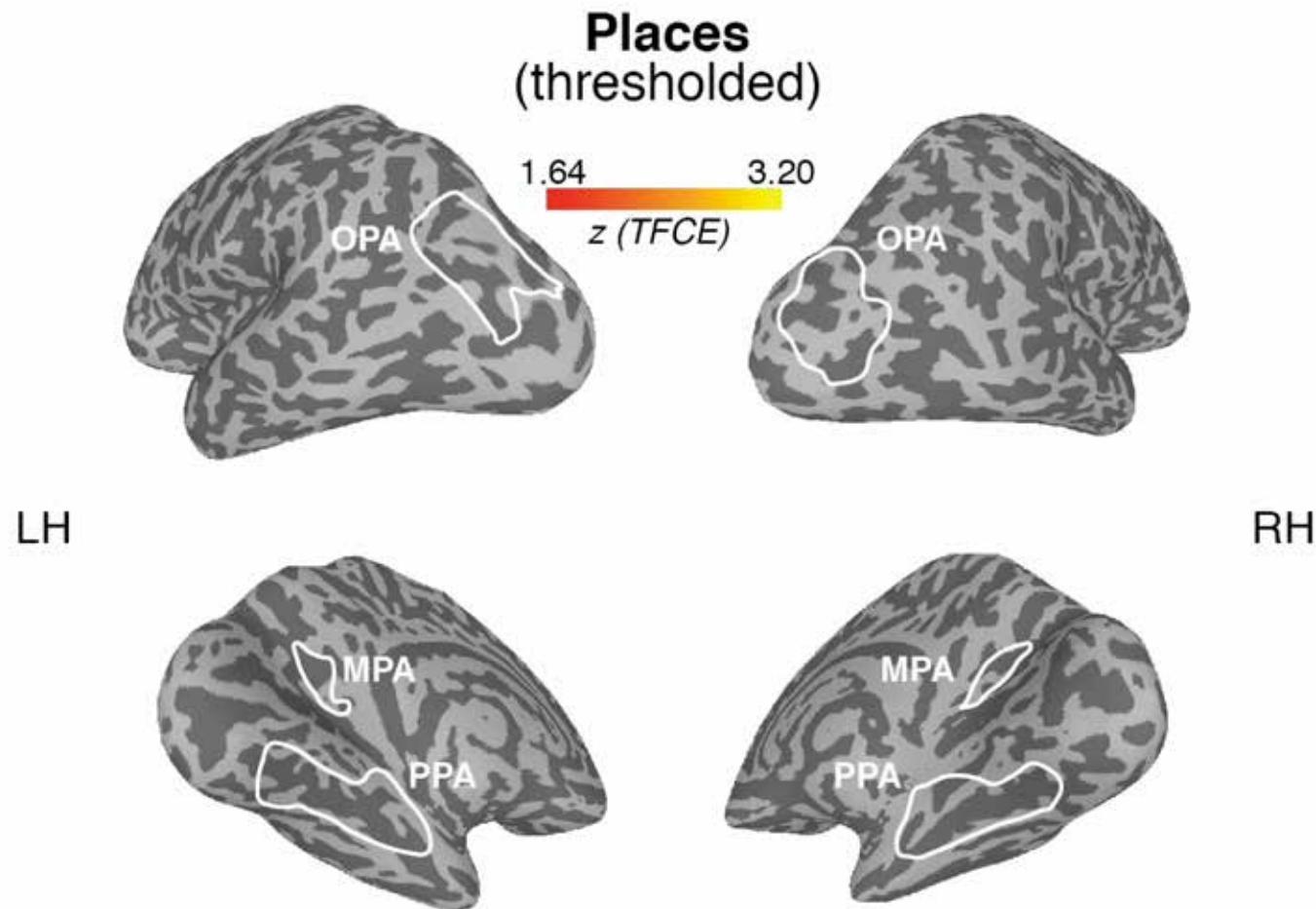
Object-trained
(ReferenceNet)



Scene-trained
(Places 205)

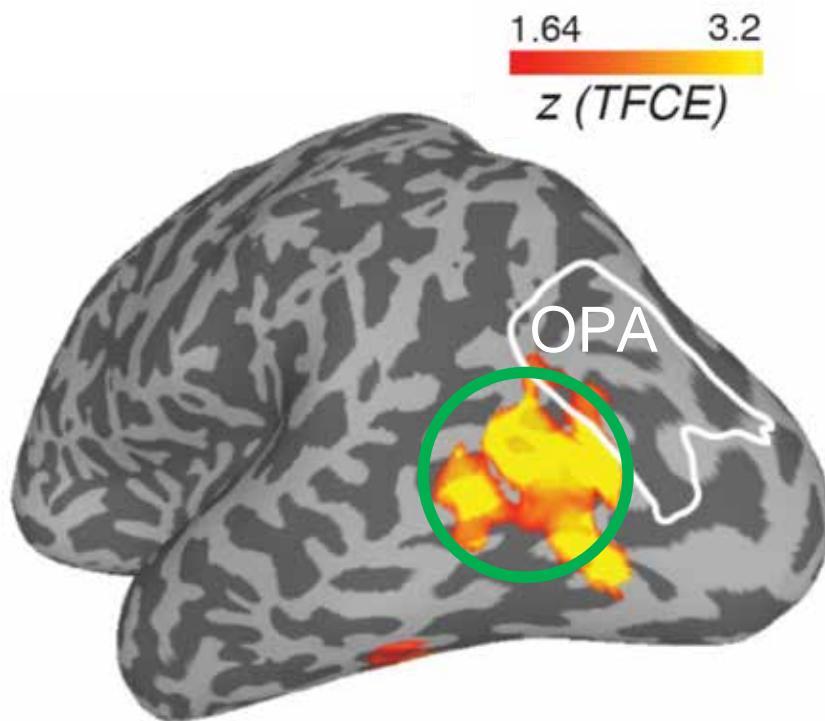


DNN correlations in scene cortex



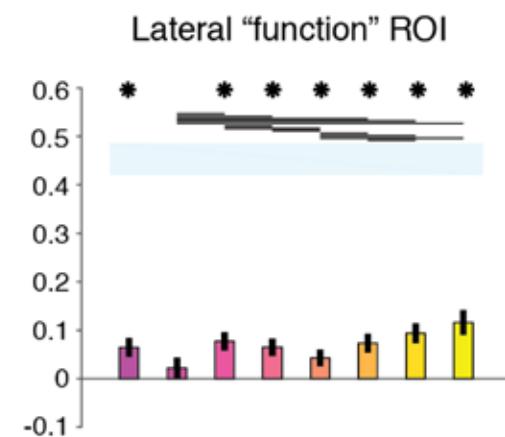
Where are the functions?

Unique correlations with function model outside scene cortex

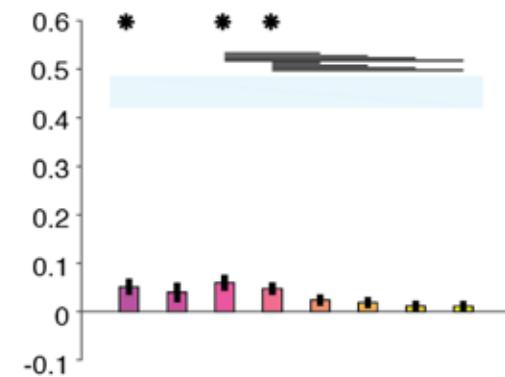


LOTC: Action perception
(Lignau & Downing, 2015)

Object-trained
(ReferenceNet)



Scene-trained
(Places 205)



Understanding scene perception

Do DNNs predict human behavioral scene perception?

Yes, but hand-labeled functions explain additional variance

Computational models of scene content

Do DNNs predict human brain responses to scenes?

Yes, but not in all regions, and not much gain with higher layers

Scene perception behavior

Not one region, but distributed mapping?

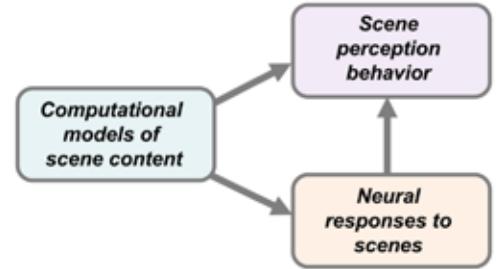
Neural responses to scenes



Outline

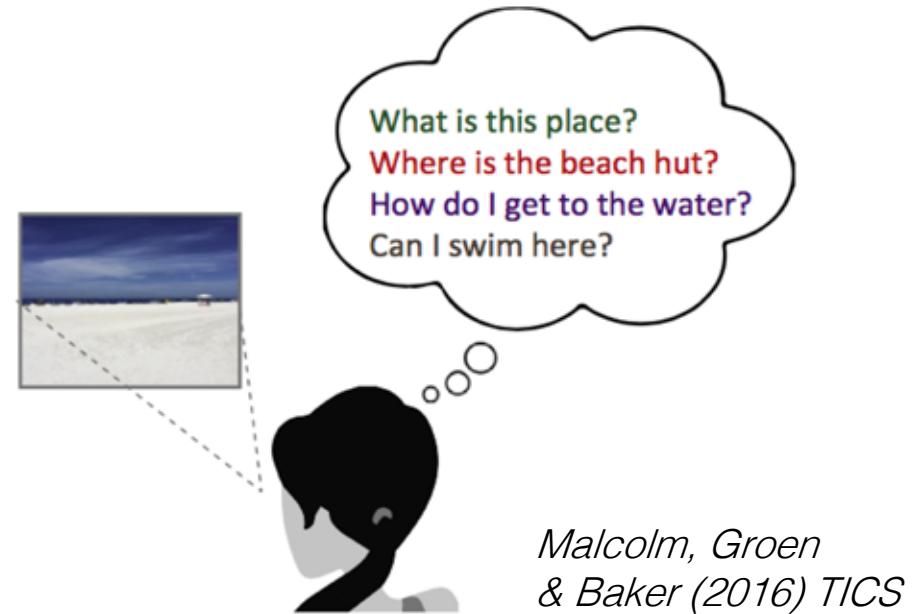
- Scene vs. object perception
- fMRI study 1: objects-in-context
- fMRI study 2: comparing multiple models
- How to move forward?

How can we close the loop?



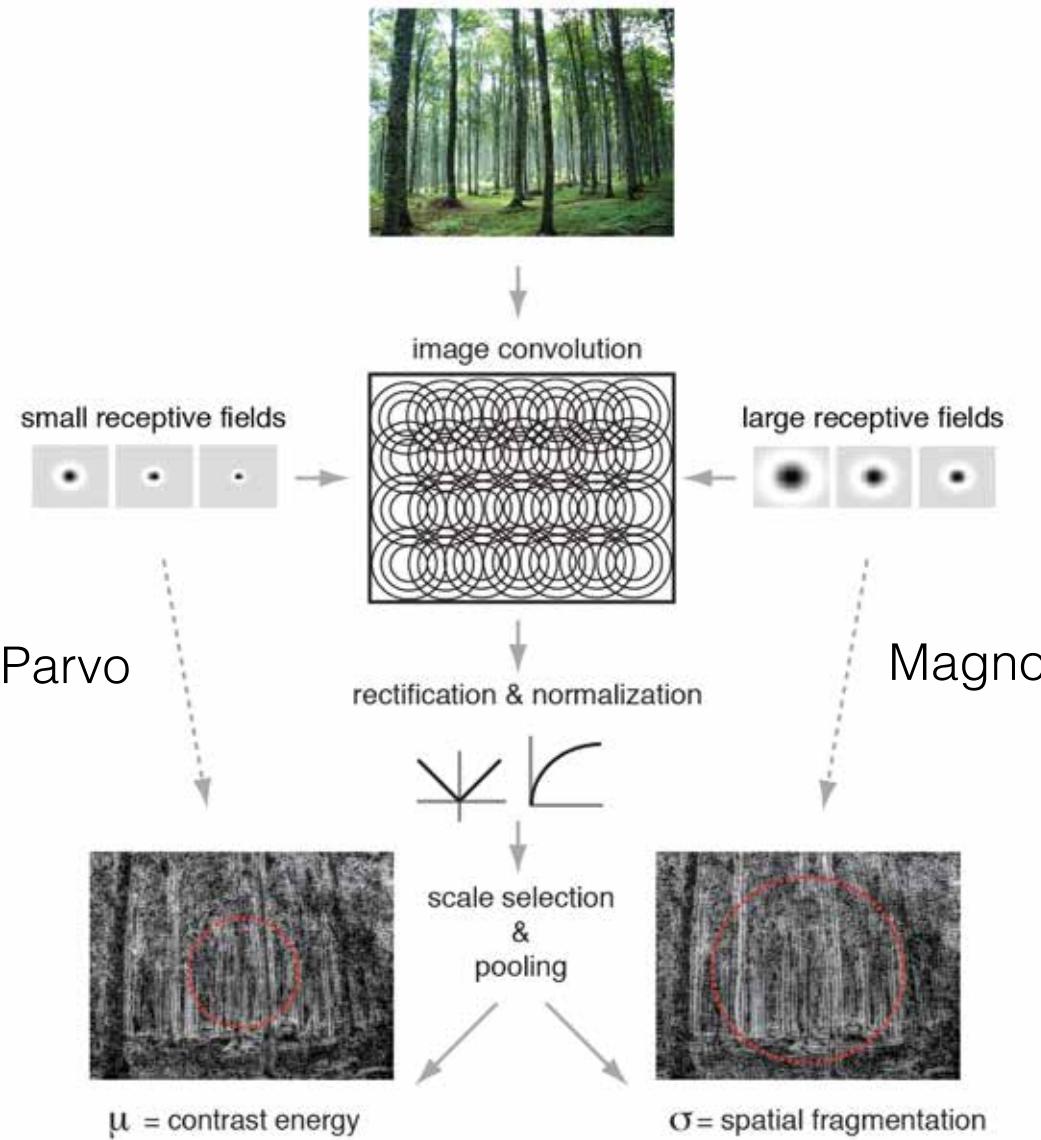
1. DNNs trained on other tasks than object or scene recognition

Recognition
Visual search
Navigation
Action

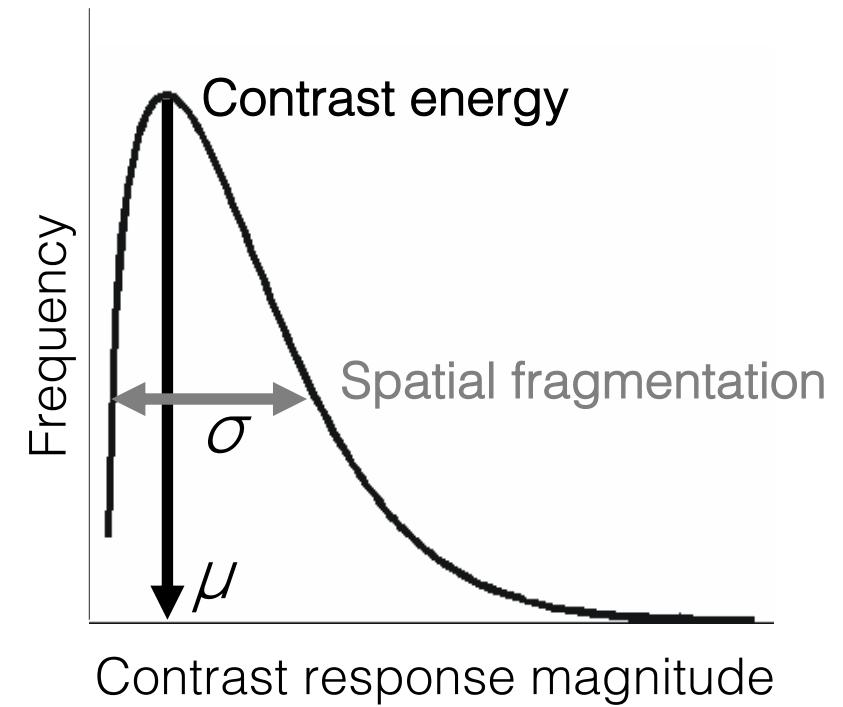


2. Other types of DNNs (e.g. RNNs), conceptual models
3. Use simpler and more explicit computational models

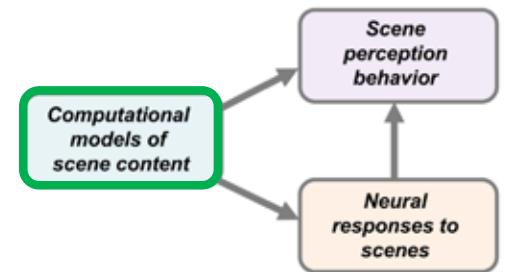
Neurophysiologically plausible model of scene statistics



Population response
to local contrast

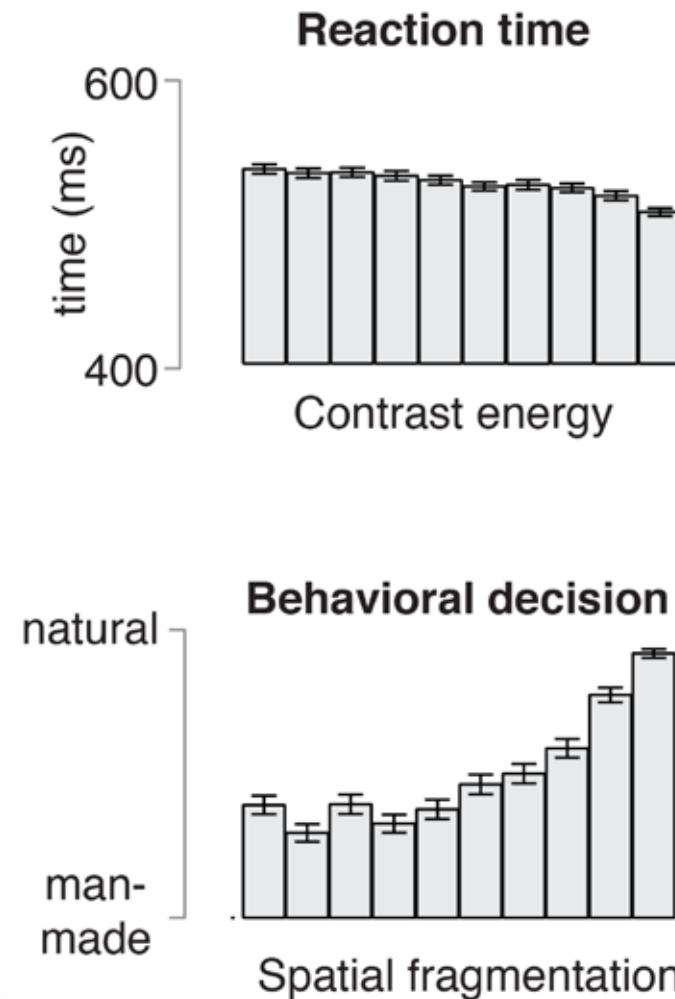
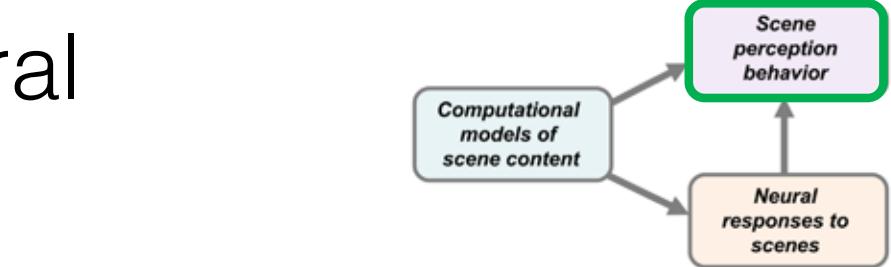
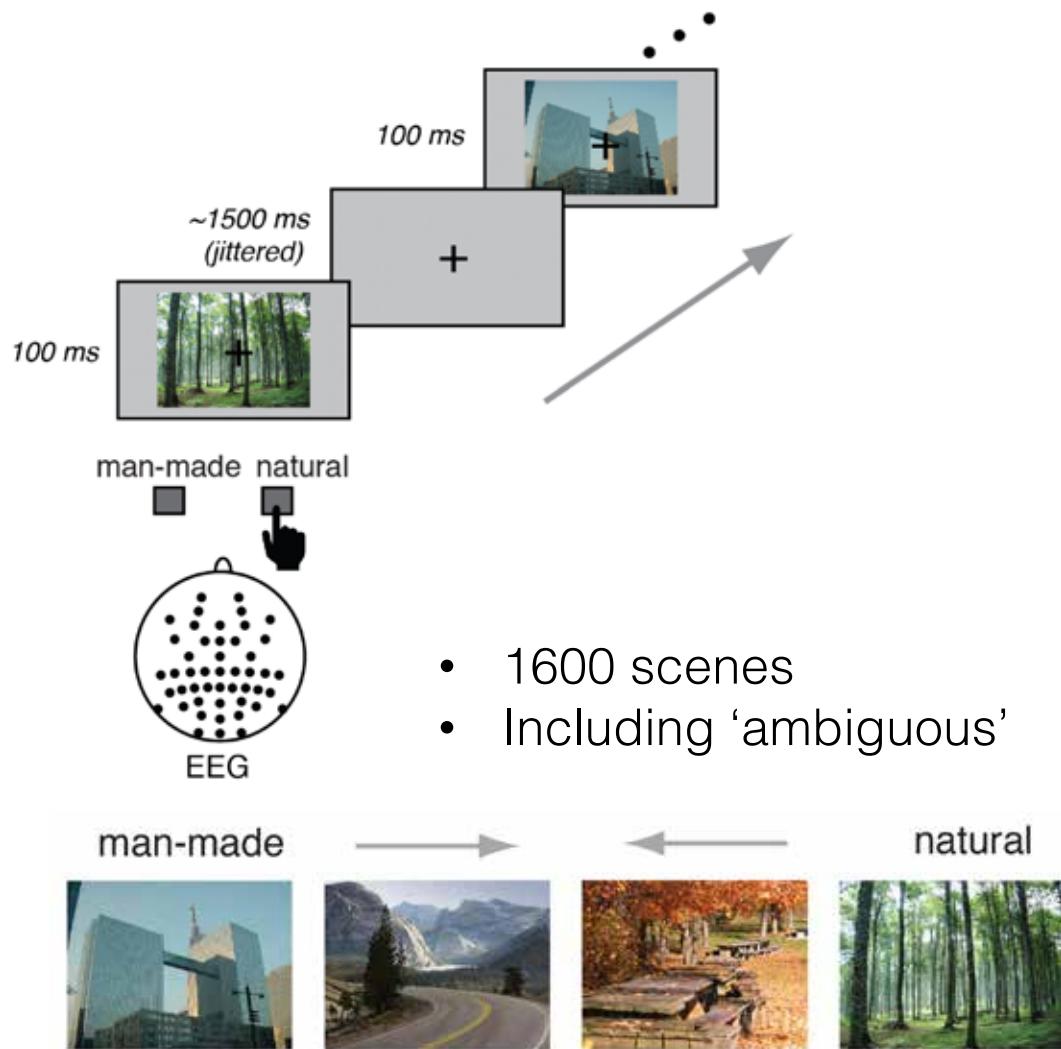


Scholte et al., (2009);
Groen et al., (2012a/b, 2013, 2017)



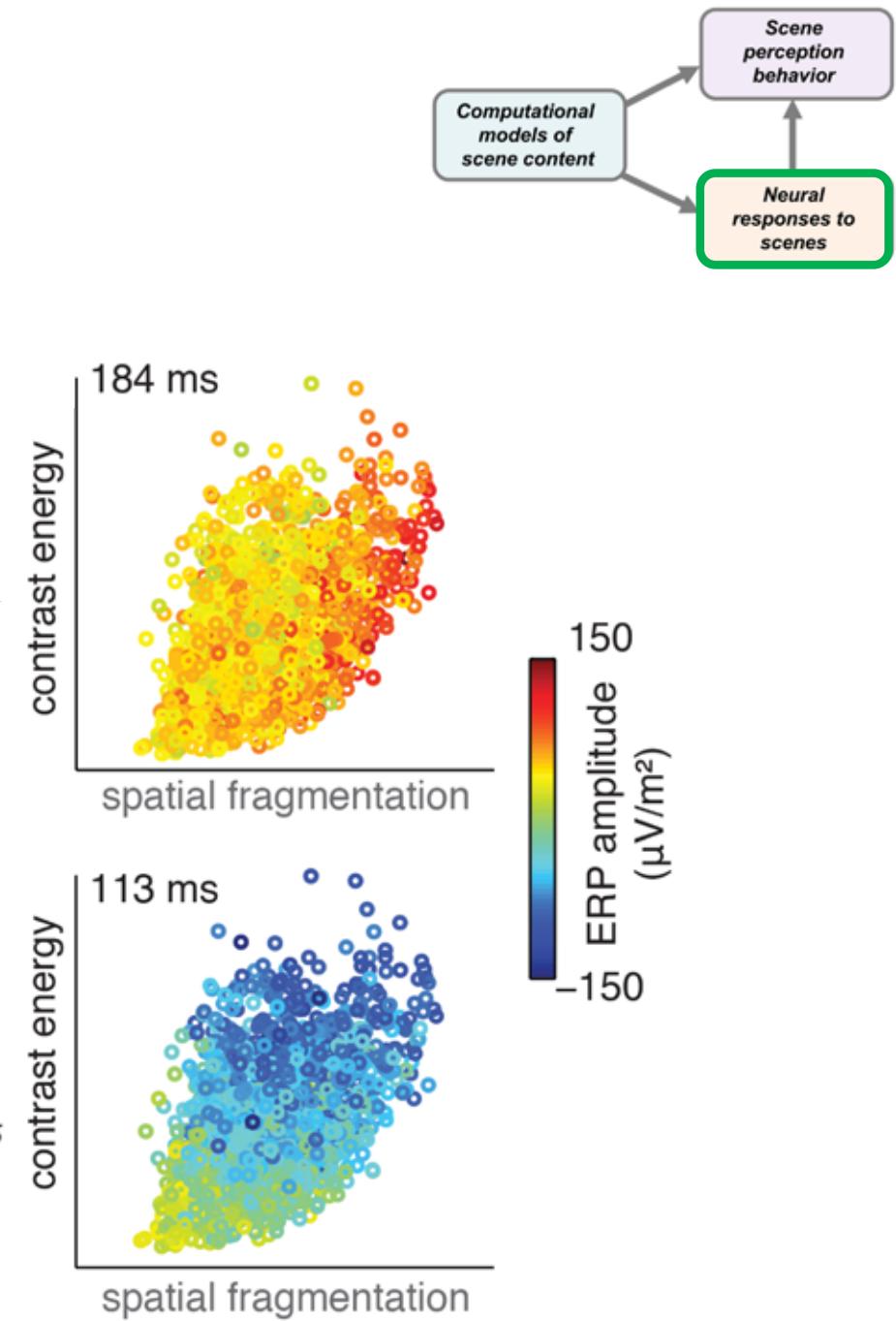
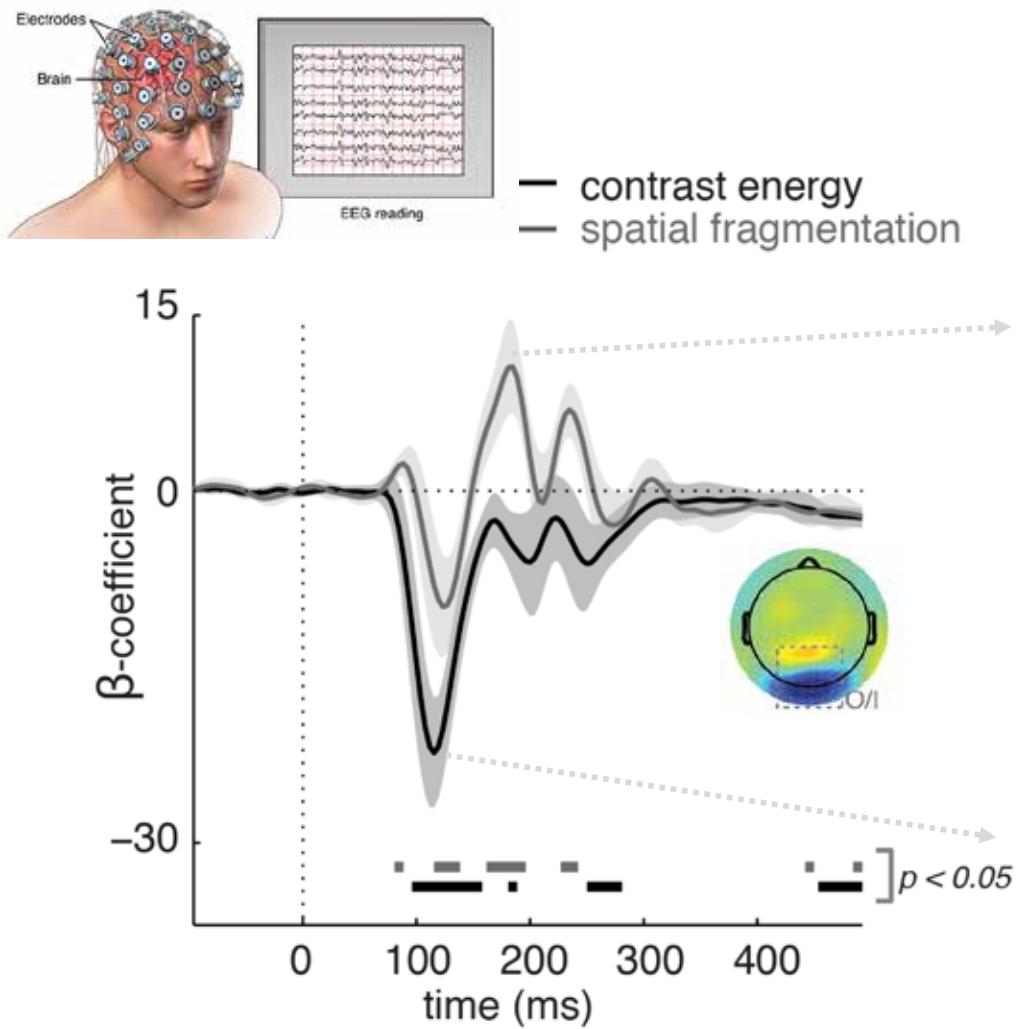
Model predicts behavioral categorization

Man-made/natural categorization task

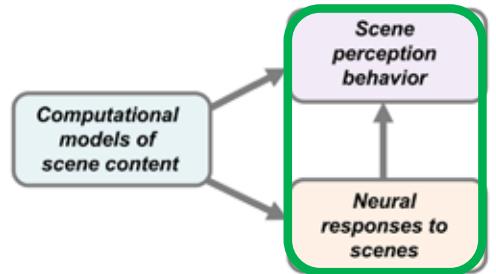


Groen et al., (2013) J Neurosci

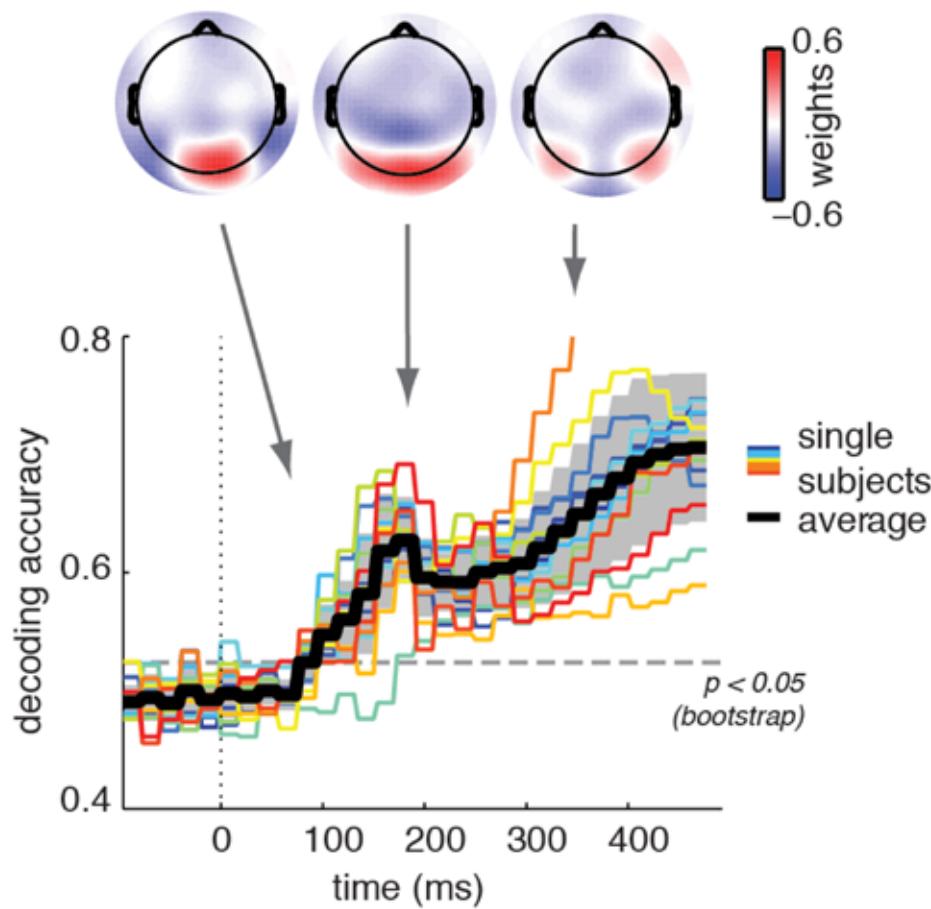
Model predicts EEG responses...



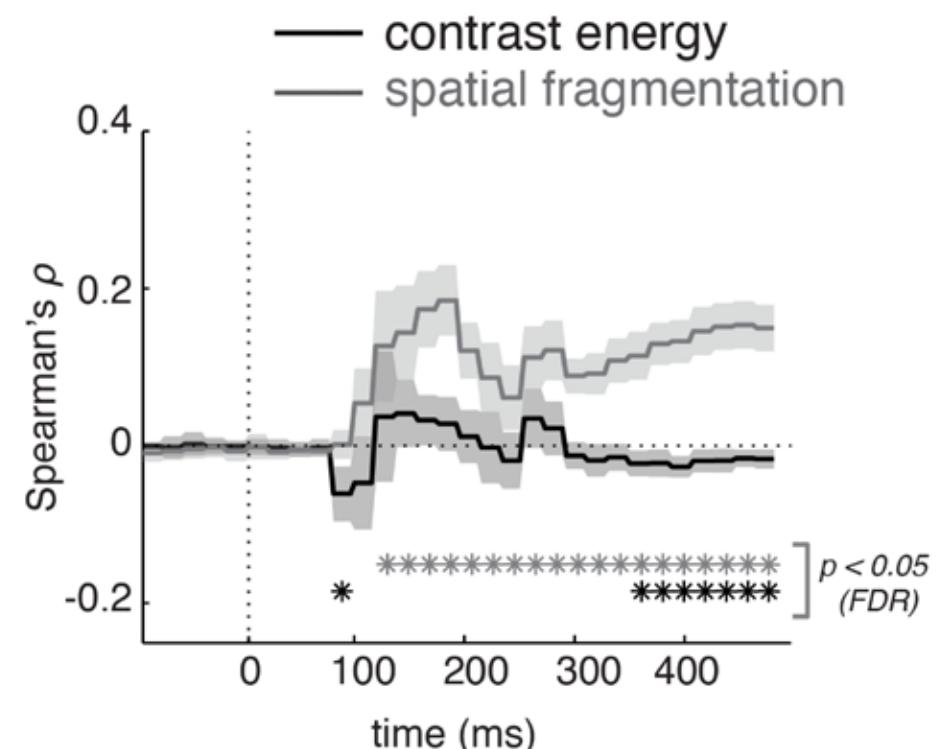
.. and decoding of behavior from EEG



Decoding accuracy

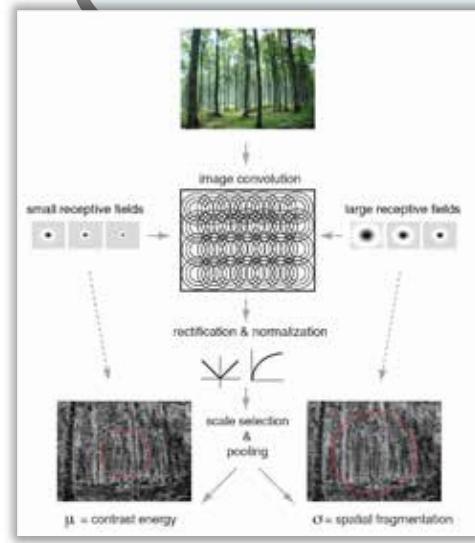


Correlation with discriminating component

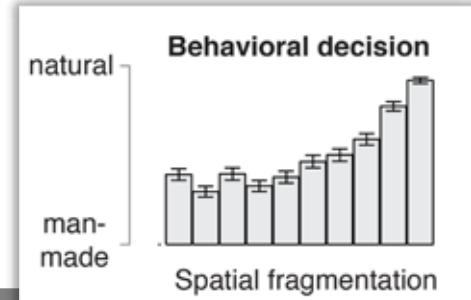


“Closed loop”

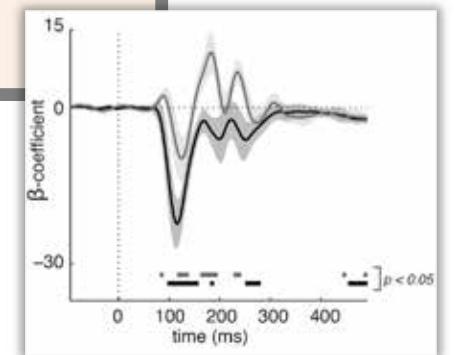
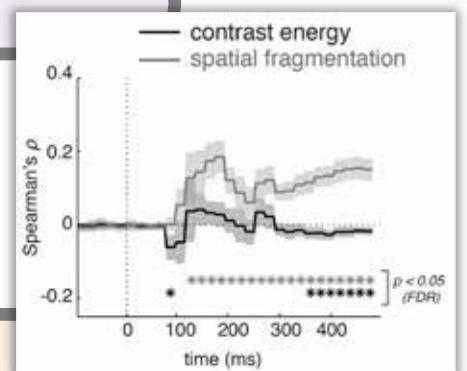
***Computational
models of
scene content***



***Scene
perception
behavior***

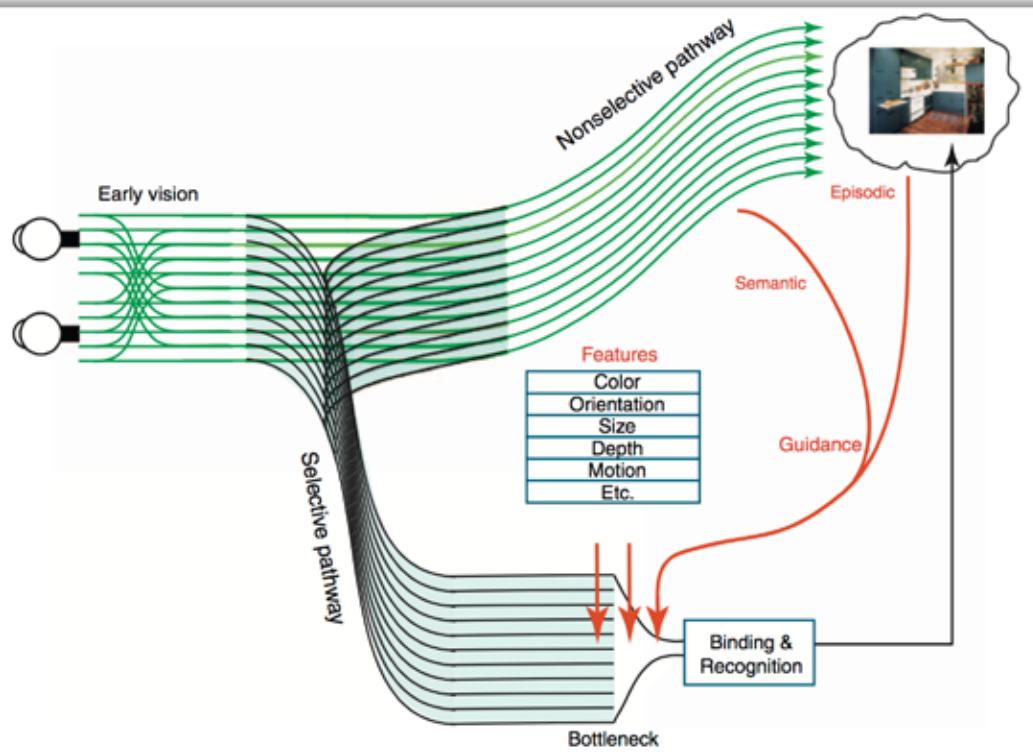


***Neural
responses to
scenes***



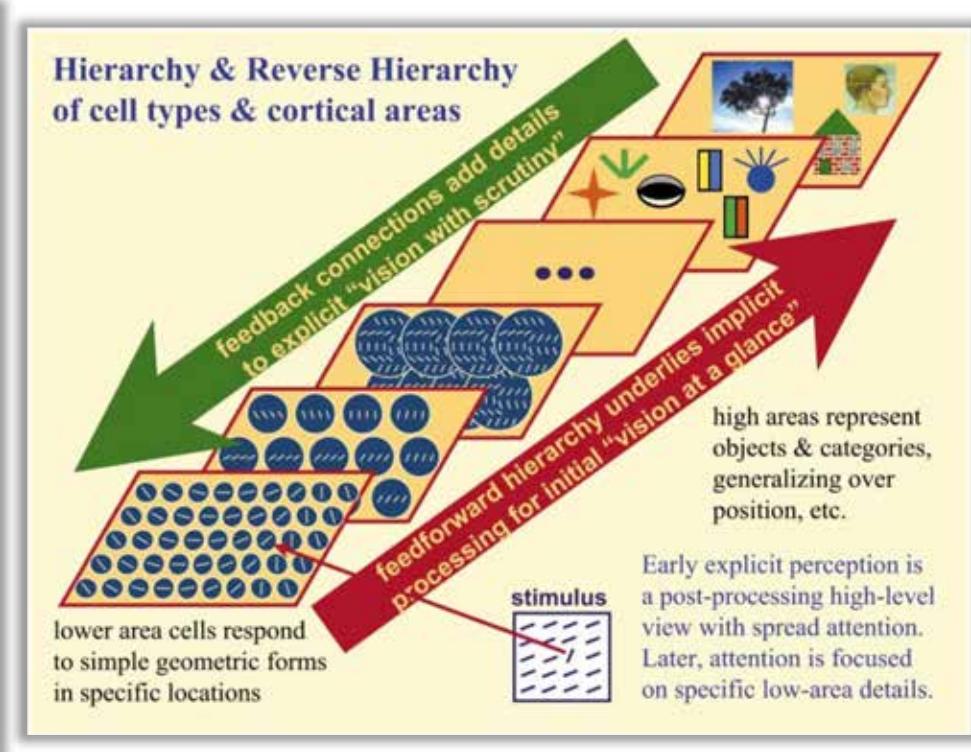
Theories of scene processing

Global vs. local



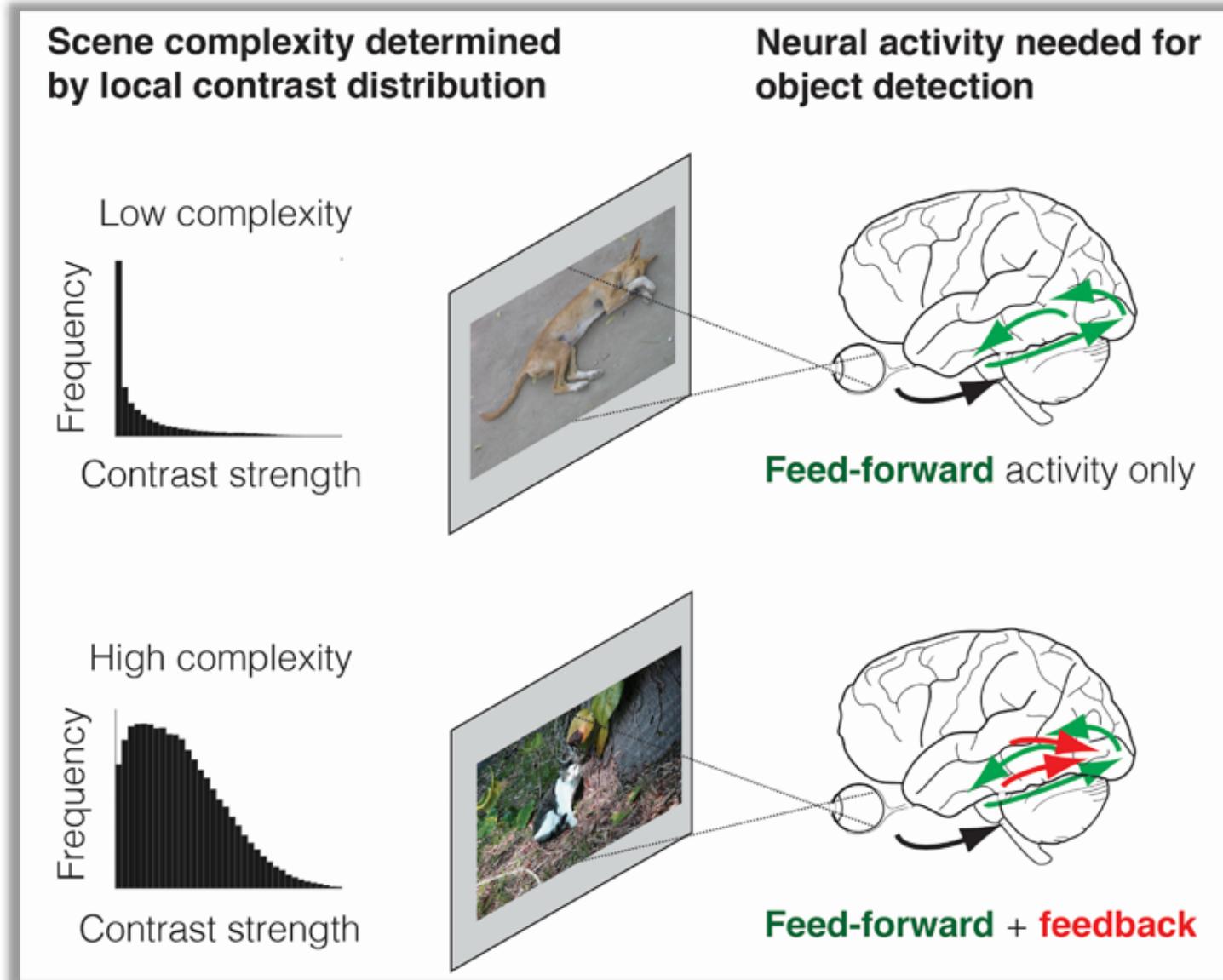
Wolfe et al., 2011

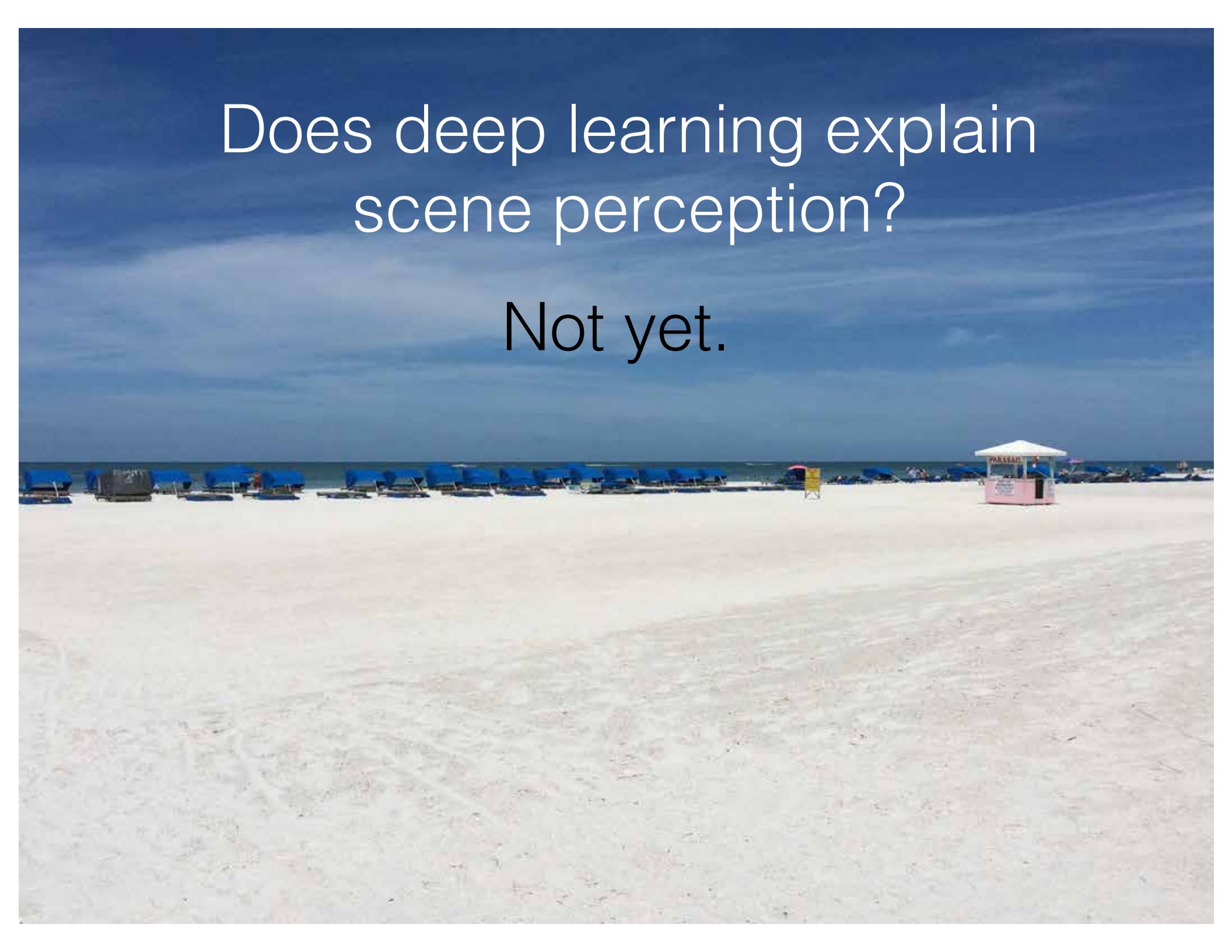
Coarse-to-fine



Hochstein & Ahissar 2002

Dynamic feedback for object recognition





Does deep learning explain
scene perception?

Not yet.

Why not?

- Scene perception entails more than object recognition
- High-level information not captured by object-trained DNNs (e.g. functions), is necessary to fully account for scene perception behavior
- Low-level information captured by image statistics, or early DNN layers, may be important for global scene representation
- DNNs capacity to predict (fMRI) responses in vTC does not automatically extend to predicting human scene perception

Thank you!

Jon Winawer
Noah Benson
Dora Hermes
Stephanie Montenegro
Jing-Yang Zhou
Chris Baker
Brett Bankson
Martin Hebart
Dwight Kravitz
Marcie King
Edward Silson
Victor Lamme
Steven Scholte
Sennay Ghebreab
Sara Jahfari
Noor Seijdel

NYU

NIH

UvA

Chris Baldassano (Columbia U)
Morgan Barense (U Toronto)
Diane Beck (U Illinois)
Tessa Dekker (UCL)
Michelle Greene (Bates College)
Assaf Harel (Wright University)
Fei-Fei Li (Stanford University)
Natalia Petridou & BAIR team (UMC Utrecht)
Kandan Ramakrishnan (MIT)



www.irisgroen.com