# Deep learning models reveal diverse computations in the retina under natural scenes

# Niru Maheswaranathan<sup>1\*</sup>, Lane McIntosh<sup>\*</sup>, David B. Kastner<sup>2</sup>, Luke Brezovec, Aran Nayebi, Surya Ganguli and Stephen A. Baccus Stanford University. <sup>1</sup>Currently at: Google Brain. <sup>2</sup>Currently at: UCSF

# Introduction

### Synthetic stimuli for probing sensory systems:

#### Advantages

Control stimulus properties parametrically (e.g. orientation) Isolate particular mechanisms/effects

### Disadvantages

May push the circuit outside its natural operating regime Often requires clever insight and design

### Goal: Understanding sensory responses to natural stimuli

- Are circuit computations discovered using structured stimuli relevant for natural vision?
- Do models trained on synthetic or natural stimuli exhibit different phenomenology?



#### **Our approach**:

- Train a deep neural network to model retinal responses to natural stimuli
- See if the model exhibits known retinal phenomenology

## CNNs accurately describe natural scene responses

• Recorded salamander ganglion cell responses to white noise and natural scene stimuli. • Fit CNN models to predict retinal spikes (blue), compared to LN models as a baseline (orange)





Population summary across *n*=37 cells:







Generalization to diverse phenomena



Remarkable match between retinal properties & internal model components Generalization beyond training distribution

- Diverse known nonlinear computations are engaged by natural stimuli Fitting models to natural stimuli is sufficient for capturing this phenomenology



# Comparing model units to retinal interneurons Example 1<sup>st</sup> layer subunits Example bipolar cell Filtered Input Example 2<sup>nd</sup> layer subunits Example adapting amacrine cells -2 -1 0 1 -0.5 0.0 0.5 Delay (s) Example model 2<sup>nd</sup> layer units Example model 1<sup>st</sup> layer units transient ON Correlations between retinal interneurons and model units Amacrine 0.4 Subunit of CNN trained on RGCs from second retina - -0.4 mm 20 40 60 80 . \_ \_ \_ \_ \_ \_ \_ \_ \_ \_ \_ \_ Amount of data (seconds) La La Time (s)

# Diversity in instantaneous receptive fields

#### Instantaneous receptive field (RF) Gradient of model neuron wrt. stimulus

I N model

- Shows the most effective stimulus at each instant of time
- Visualizes stimulus-specific sensitivity

# Conclusions

- Implications for using deep networks to study neural systems:
- Implications for retinal encoding of natural scenes: