# Strong and Precise Modulation of Human Percepts via Robustified ANNs



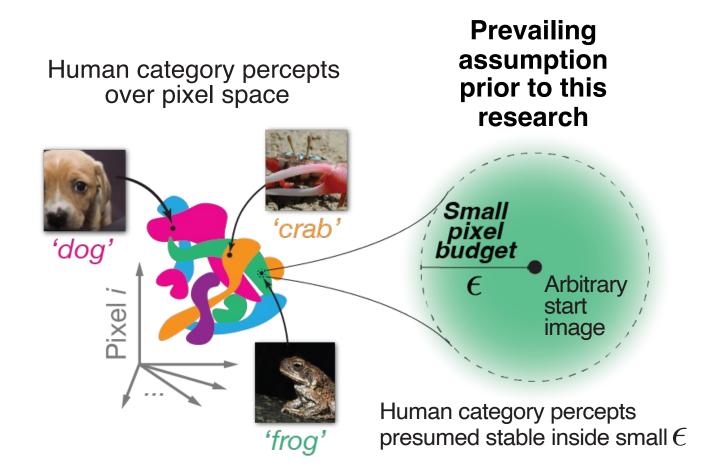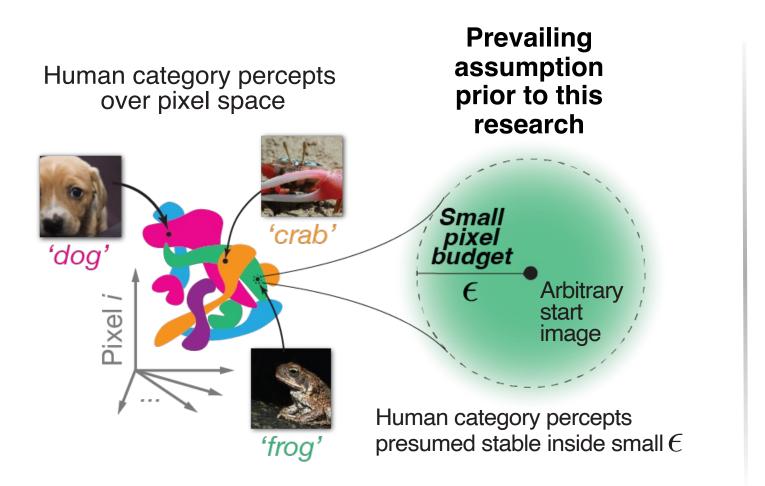**G. Gaziv***    **M. Lee***    **Jim DiCarlo**

*Dept of Brain & Cognitive Sciences, MIT*

# Is perception robust as believed?



Human category percepts over pixel space

Prevailing assumption prior to this research

'dog'

'crab'

'frog'

Pixel $i$

...

Small pixel budget $\epsilon$

Arbitrary start image

Human category percepts presumed stable inside small $\epsilon$

# The truth of human biology near any image start point



Human category percepts over pixel space

**Prevailing assumption prior to this research**

'dog'

'crab'

'frog'

Pixel *i*

*Small pixel budget*

$\epsilon$

Arbitrary start image

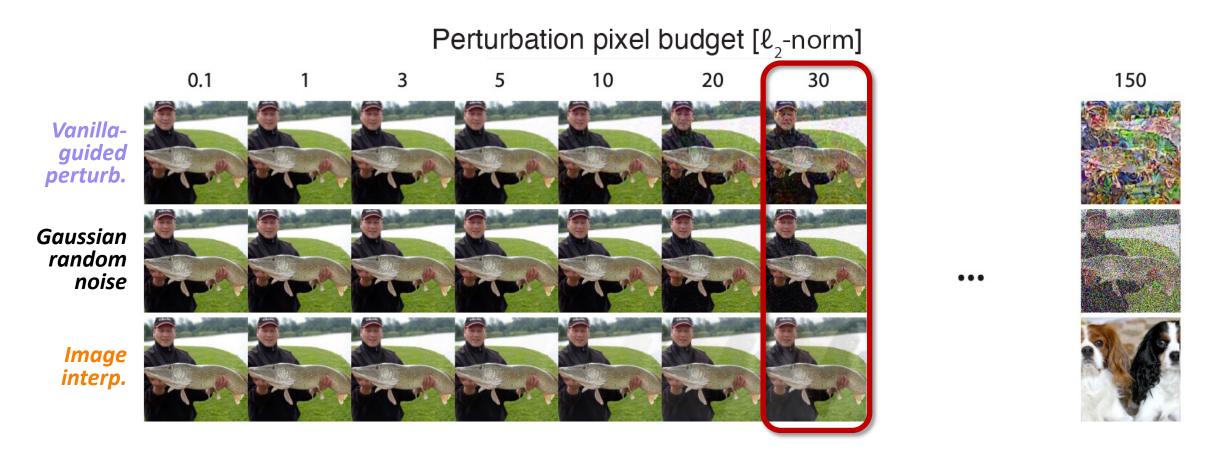Human category percepts presumed stable inside small $\epsilon$

# Robustified models allow for precise Targeted Modulation of human behavior in the "human-presumed-insensitive" pixel budget regime

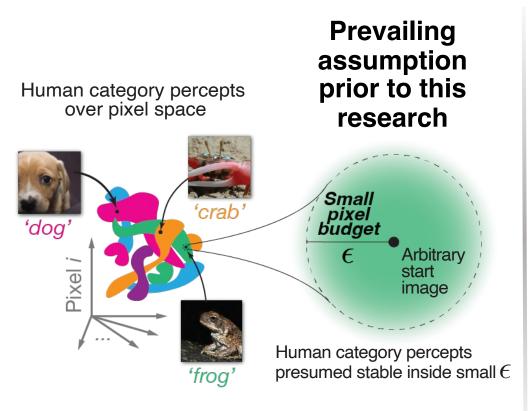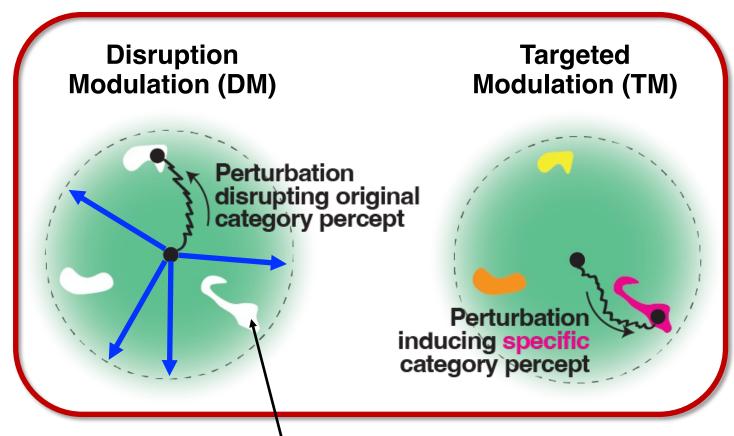# Targeted Modulation successfully modulates human reports for _arbitrary_ start image



$\epsilon = 30$

# Humans are (indeed) insensitive to random, vanilla-guided or interpolation-based perturbations in the "low norm" regime



Perturbation pixel budget [$\ell_2$-norm]

# ANNs reveal "wormholes" between human category percepts



**Human category percepts over pixel space**

'dog' 'crab' 'frog'

Pixel $i$

**Prevailing assumption prior to this research**

*Small pixel budget* $\epsilon$

Arbitrary start image

Human category percepts presumed stable inside small $\epsilon$

**Disruption Modulation (DM)**

Perturbation disrupting original category percept

**Targeted Modulation (TM)**

Perturbation inducing **specific** category percept

Vanilla ANN vision models were unable to locate undetected "wormholes" that we now easily find

Code: https://github.com/ggaziv/Wormholes
Project Page: https://himjl.github.io/pwormholes

*Gaziv*, Lee* & DiCarlo, **NeurIPS** (2023)*