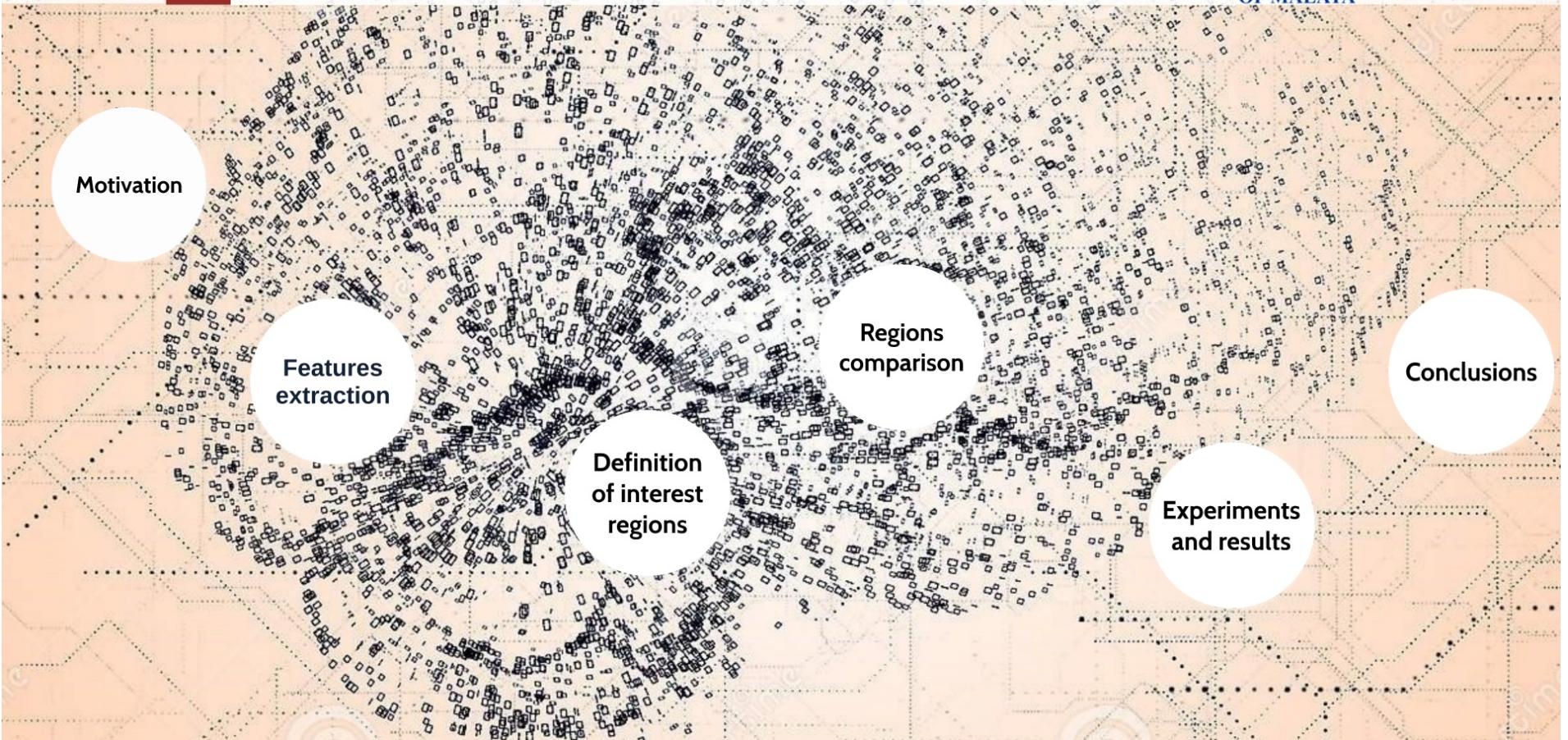
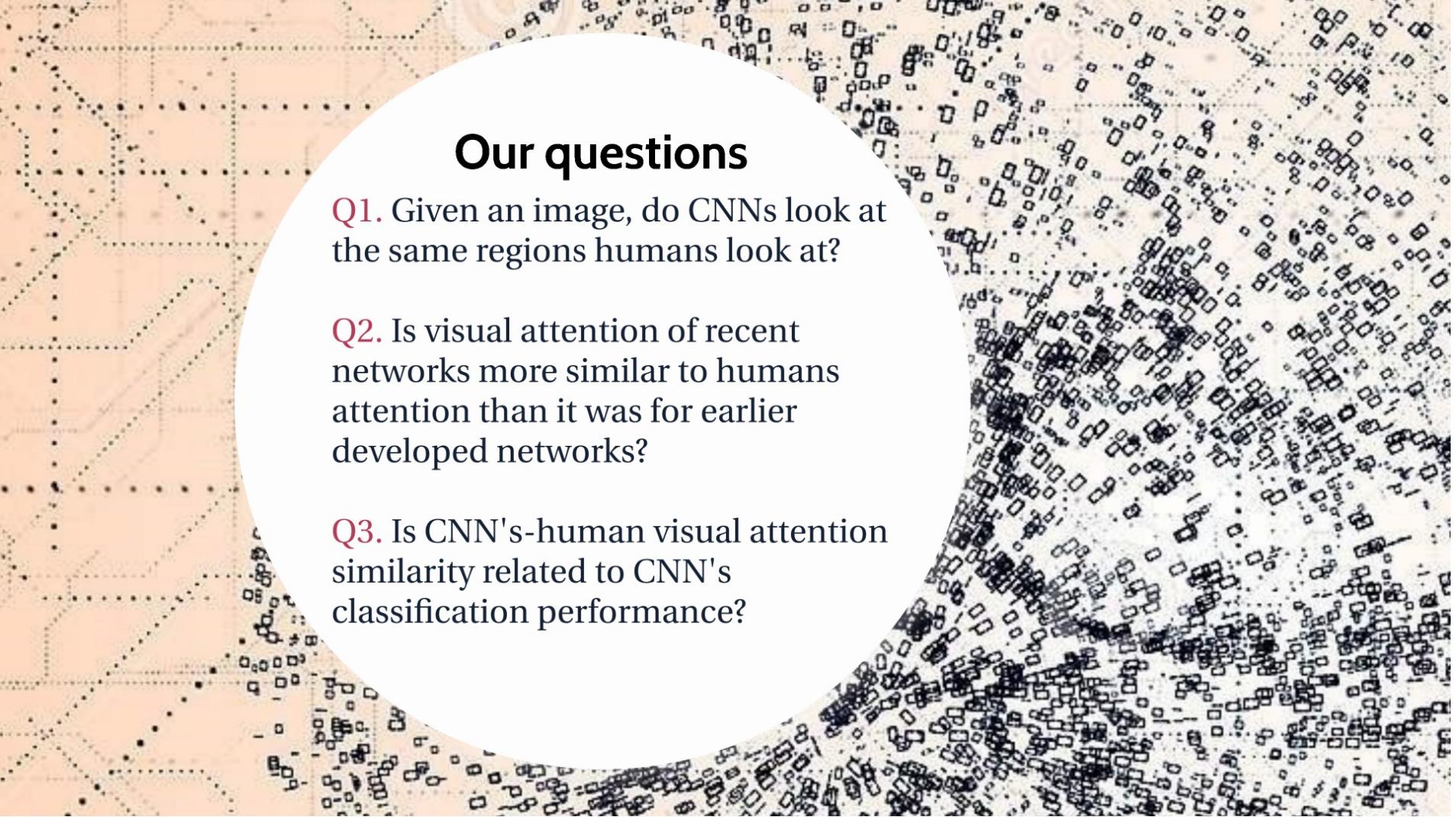


# From early biological models to CNNs: do they look where humans look?

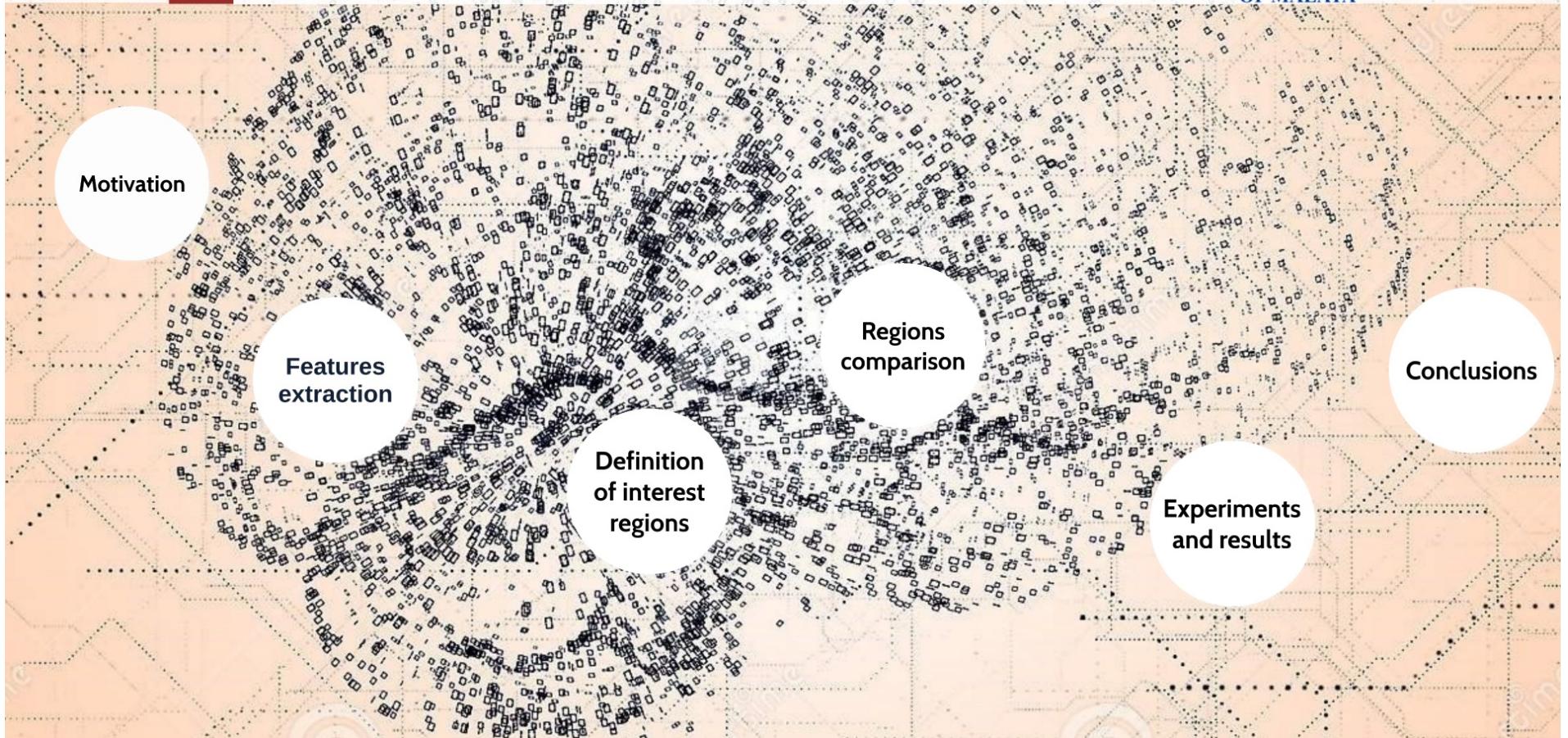


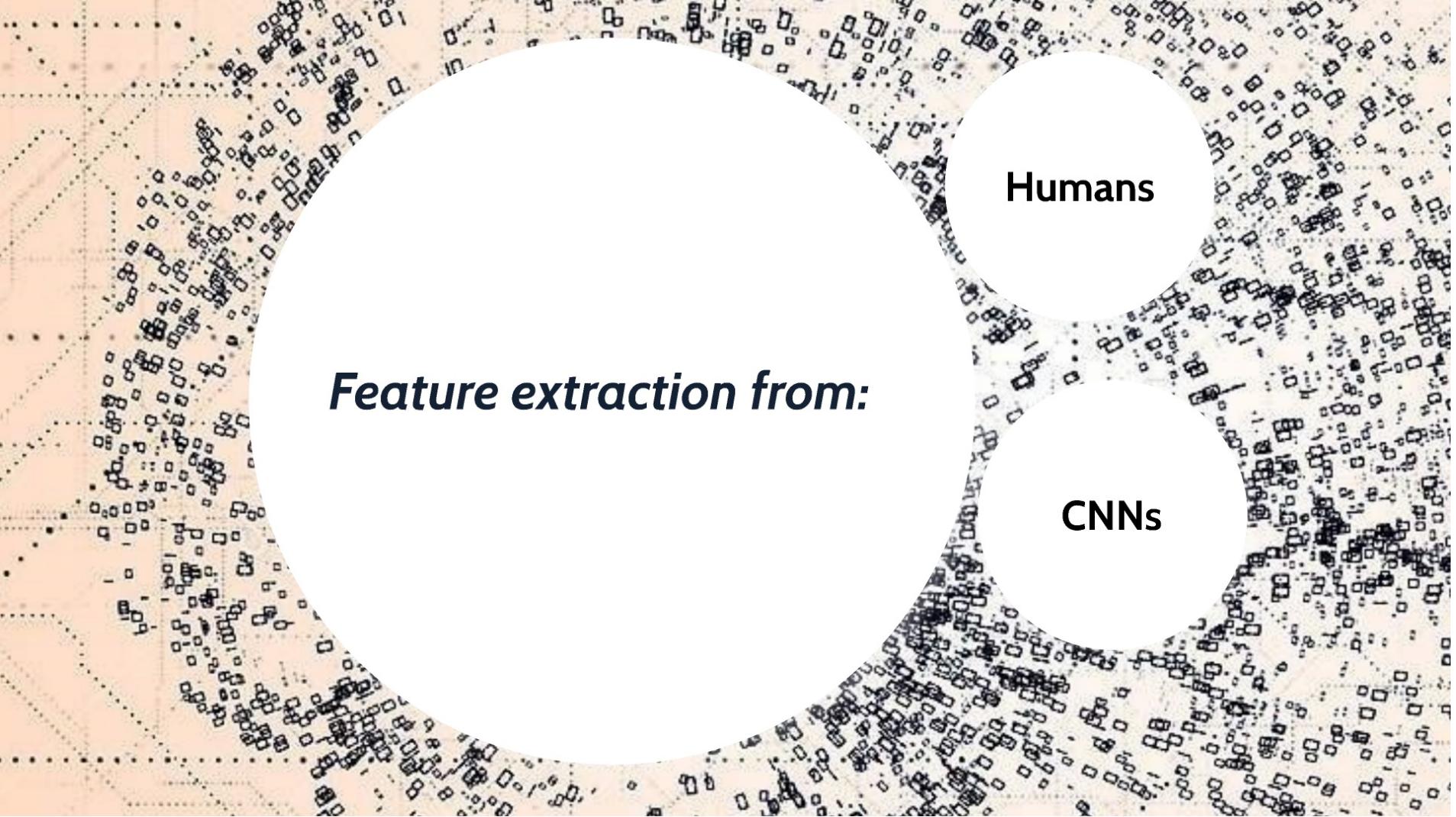


## Our questions

- Q1.** Given an image, do CNNs look at the same regions humans look at?
- Q2.** Is visual attention of recent networks more similar to humans attention than it was for earlier developed networks?
- Q3.** Is CNN's-human visual attention similarity related to CNN's classification performance?

# From early biological models to CNNs: do they look where humans look?





***Feature extraction from:***

**Humans**

**CNNs**

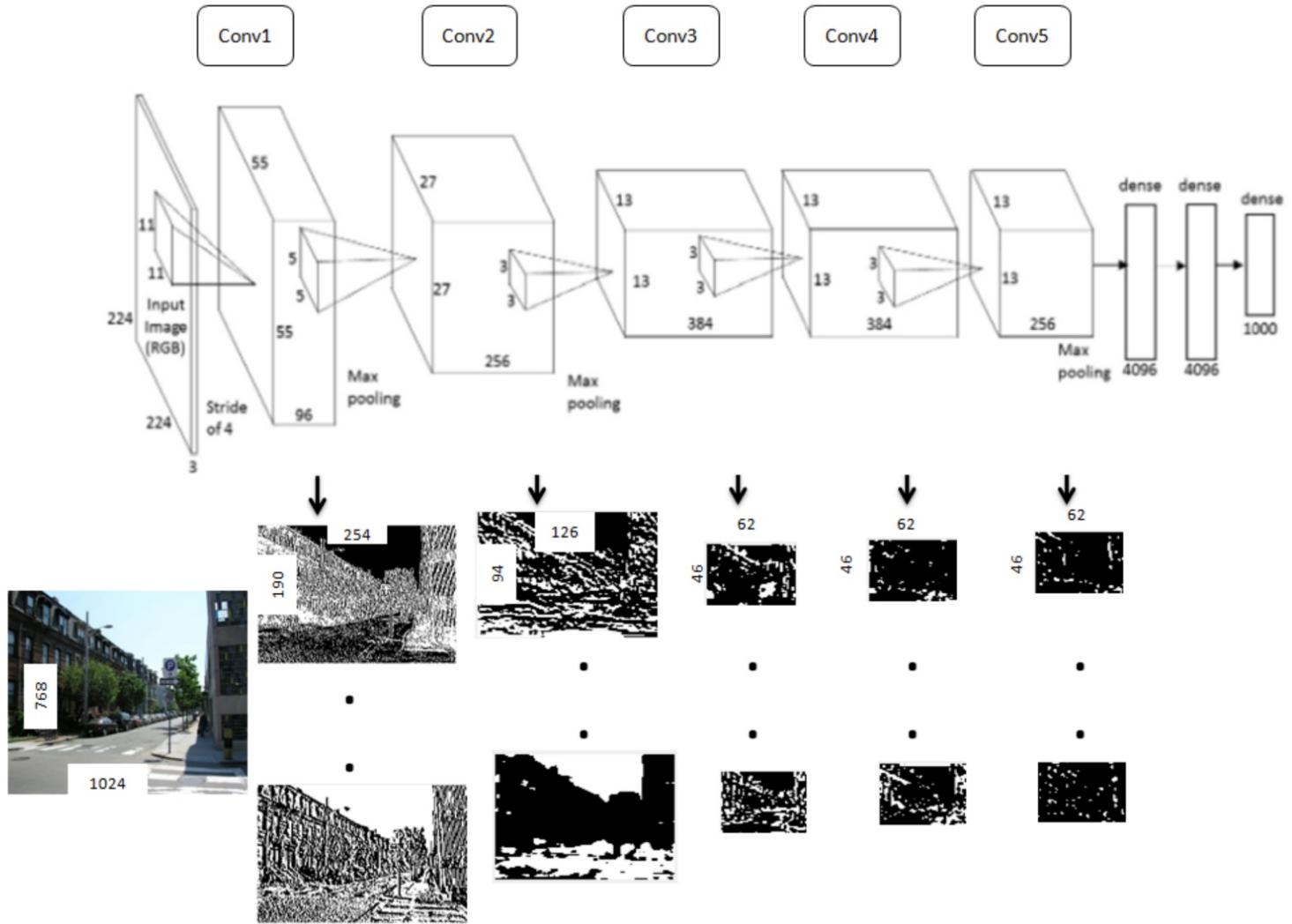
# MIT Eye Tracking Database

- Gaze tracking paths (saccades and fixations)
  - 15 users sitting in front of a screen
  - 1003 images freely viewed
  - 3 seconds per image with 1 second of gray screen in between



# AlexNet

- 8 layers
- 5 convolutional
- 3 fully connected



## Example of human fixations and CNNs feature points.

Clockwise:

- original image
- humans fixations
- HMAX features

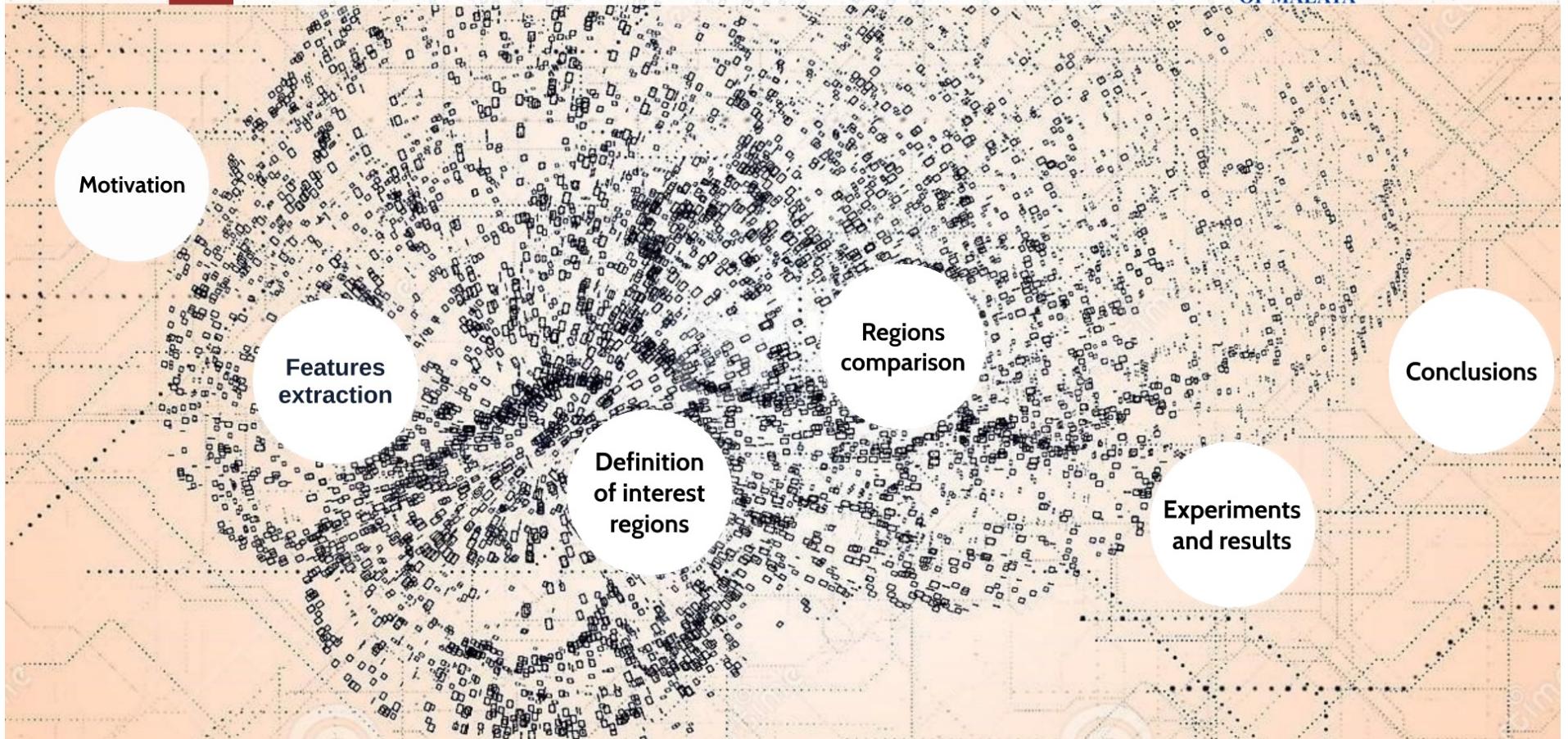
first, middle and last layers

features of:

- VGG-19
- Resnet-50
- Efficientnet.



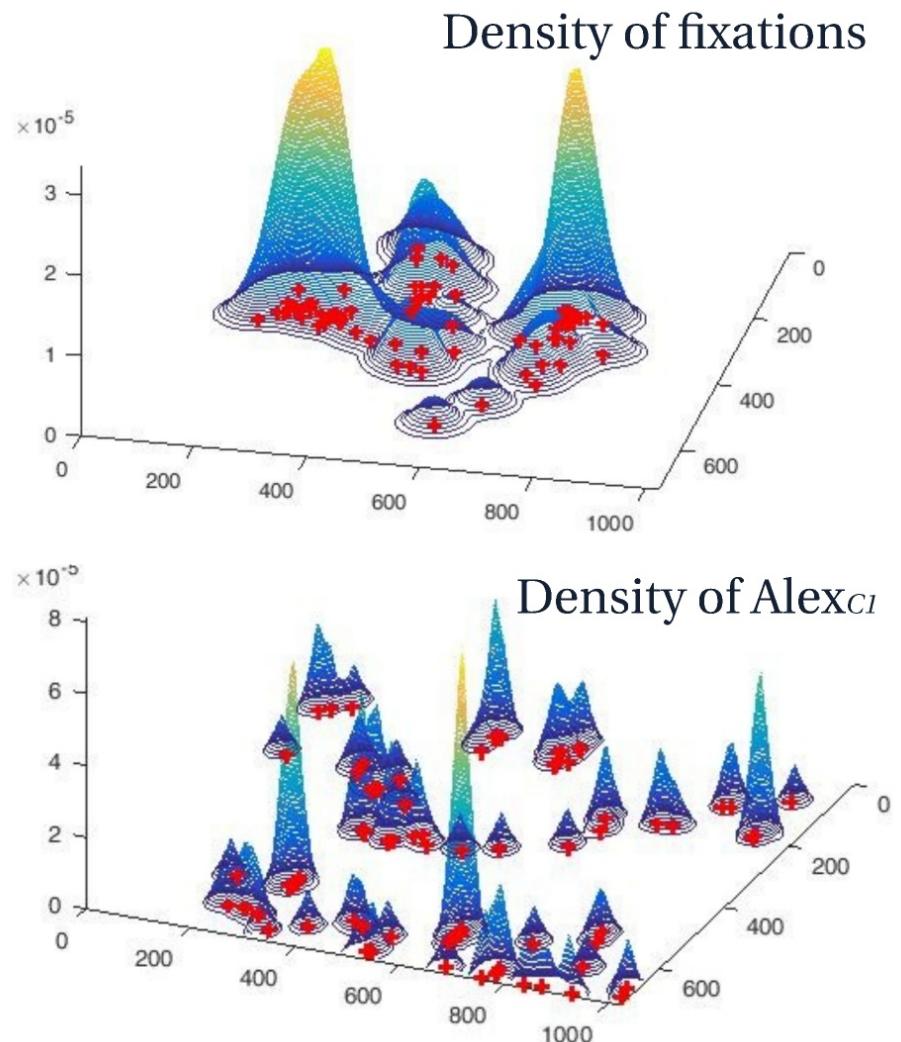
# From early biological models to CNNs: do they look where humans look?



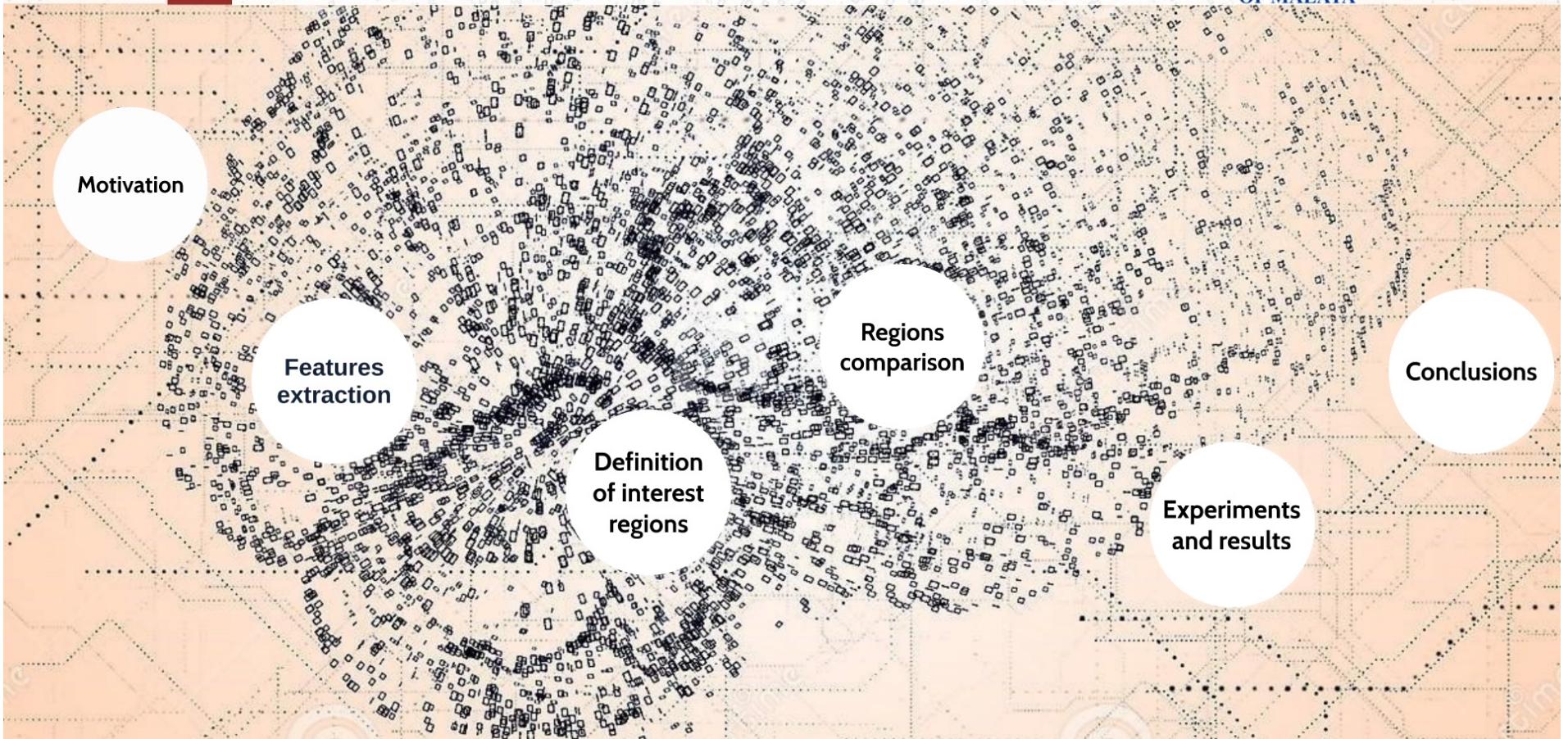
# Regions of interest

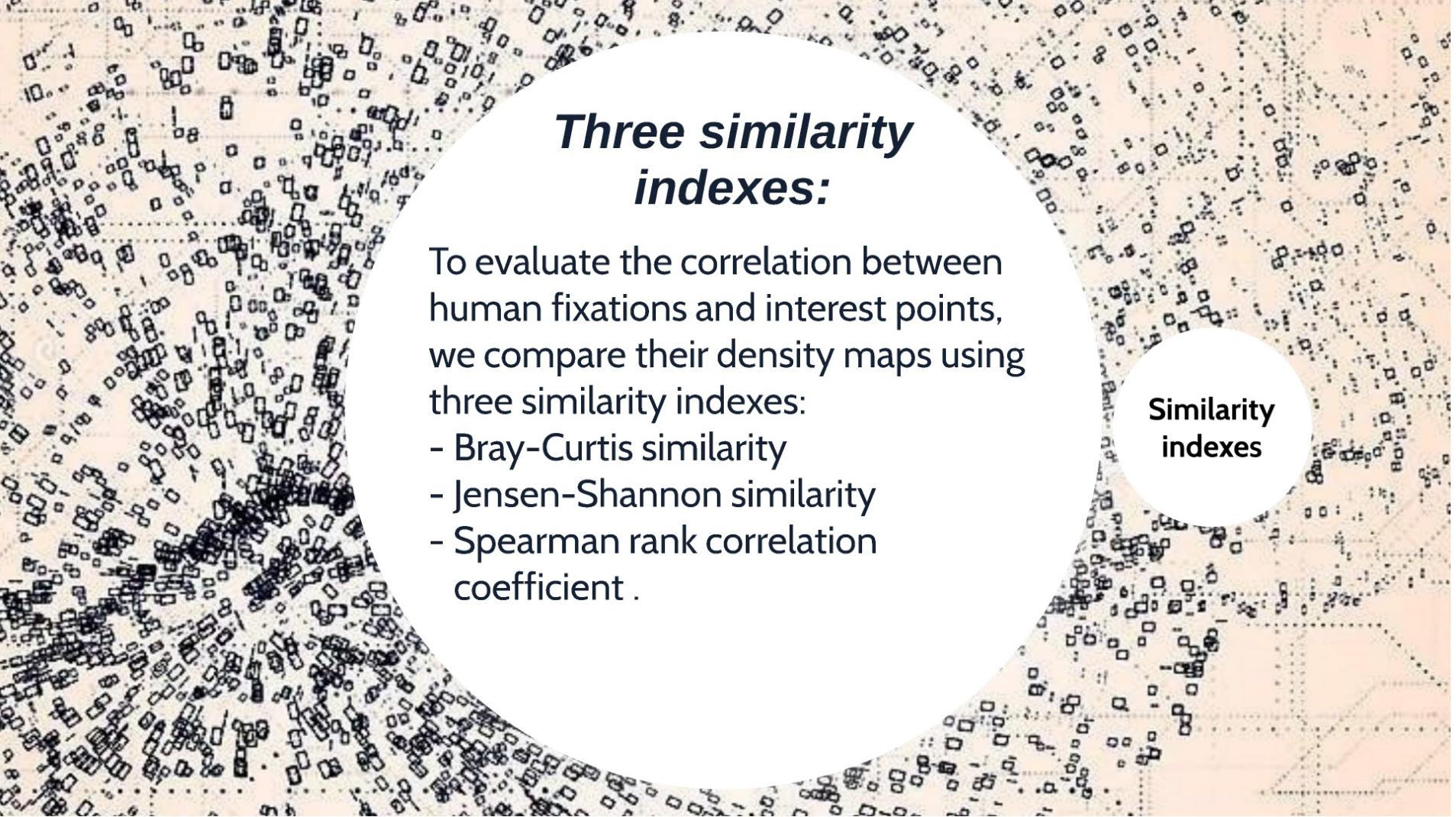
- Let  $I$  be an image and  $F(I)$  a set of features.
- The region of interest is defined to be the support of the probability density function of the distribution of the feature points  $F(I)$ .
- The density  $f_F(I)$  is estimated through a KDE (Kernel Density Estimation) based on a linear diffusion model proposed by Botev *et al.* [Kernel density estimation via diffusion, *The annals of statistics*, 2010]

Red points: Human fixations  
Blue points: AlexNet<sub>C1</sub> features



# From early biological models to CNNs: do they look where humans look?





## *Three similarity indexes:*

To evaluate the correlation between human fixations and interest points, we compare their density maps using three similarity indexes:

- Bray-Curtis similarity
- Jensen-Shannon similarity
- Spearman rank correlation coefficient .

# Similarity indexes

Bray Curtis similarity:

$$BC_{1,2} = 1 - \frac{\sum_{i=1}^n |f_1(\mathbf{x}_i) - f_2(\mathbf{x}_i)|}{\sum_{i=1}^n f_1(\mathbf{x}_i) + f_2(\mathbf{x}_i)}$$

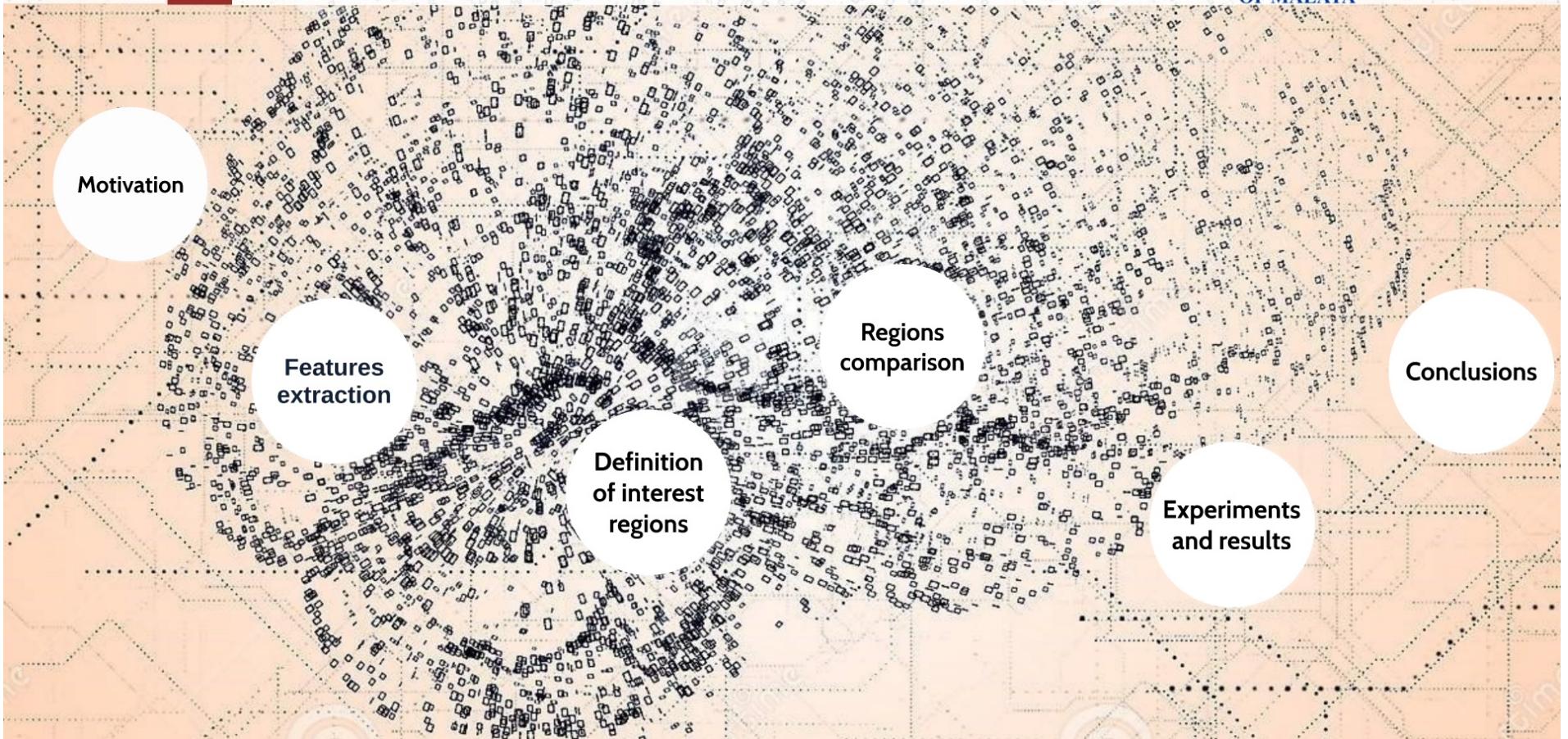
Jensen Shannon  
Divergence:

$$JSD_{1,2} = \sum_{\mathbf{x}_i} (f_1(\mathbf{x}_i) - f_2(\mathbf{x}_i)) \log \frac{f_1(\mathbf{x}_i)}{f_2(\mathbf{x}_i)}$$

Jensen Shannon  
similarity:

$$JS_{1,2} = 1 - \bar{JSD}_{1,2}$$

## From early biological models to CNNs: do they look where humans look?

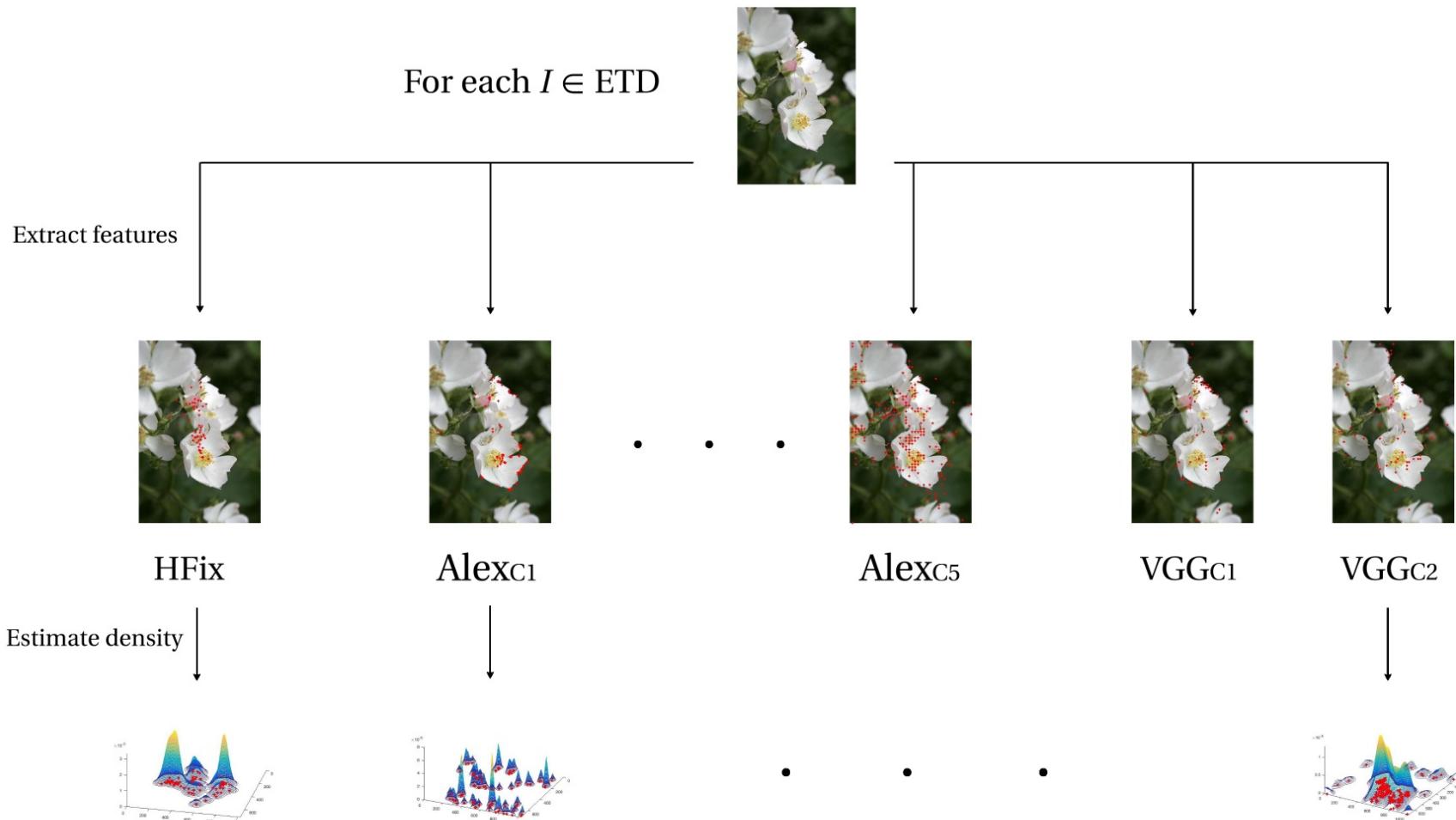


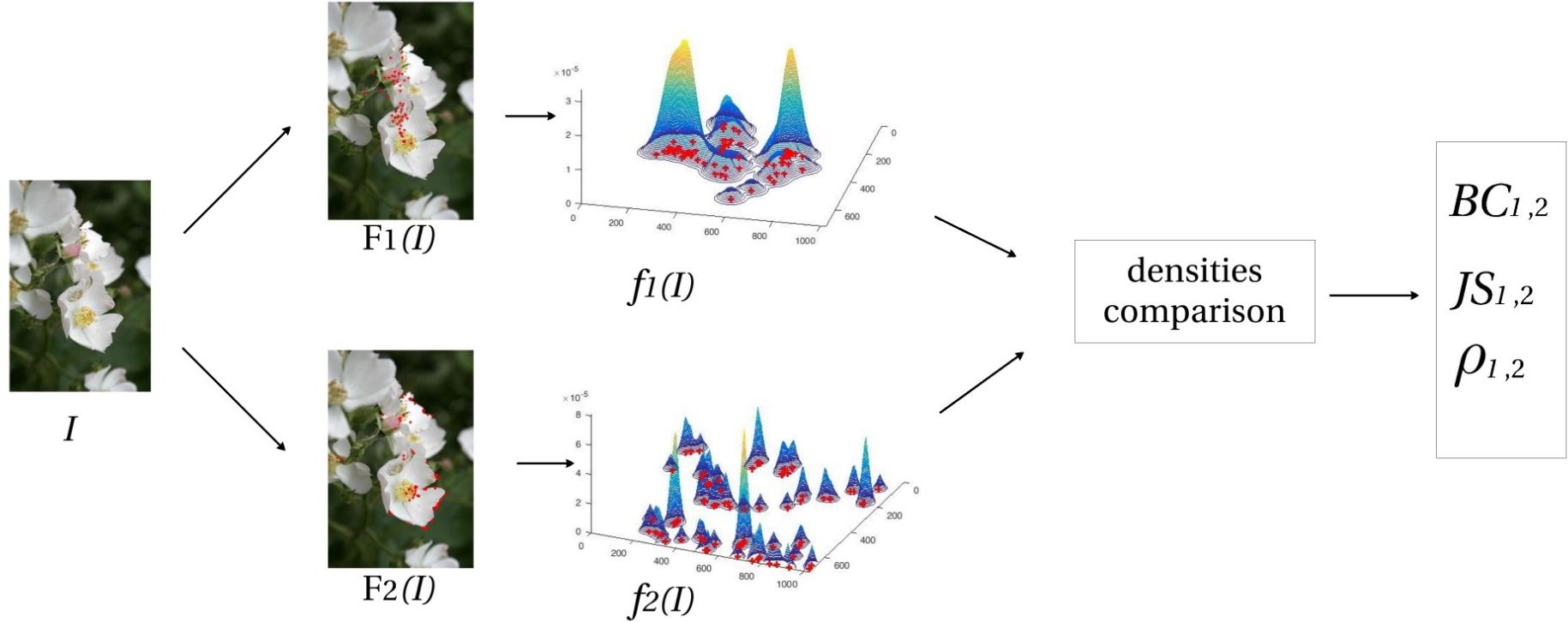


# **Experiments and results**

**Experimental  
protocol**

**Results**

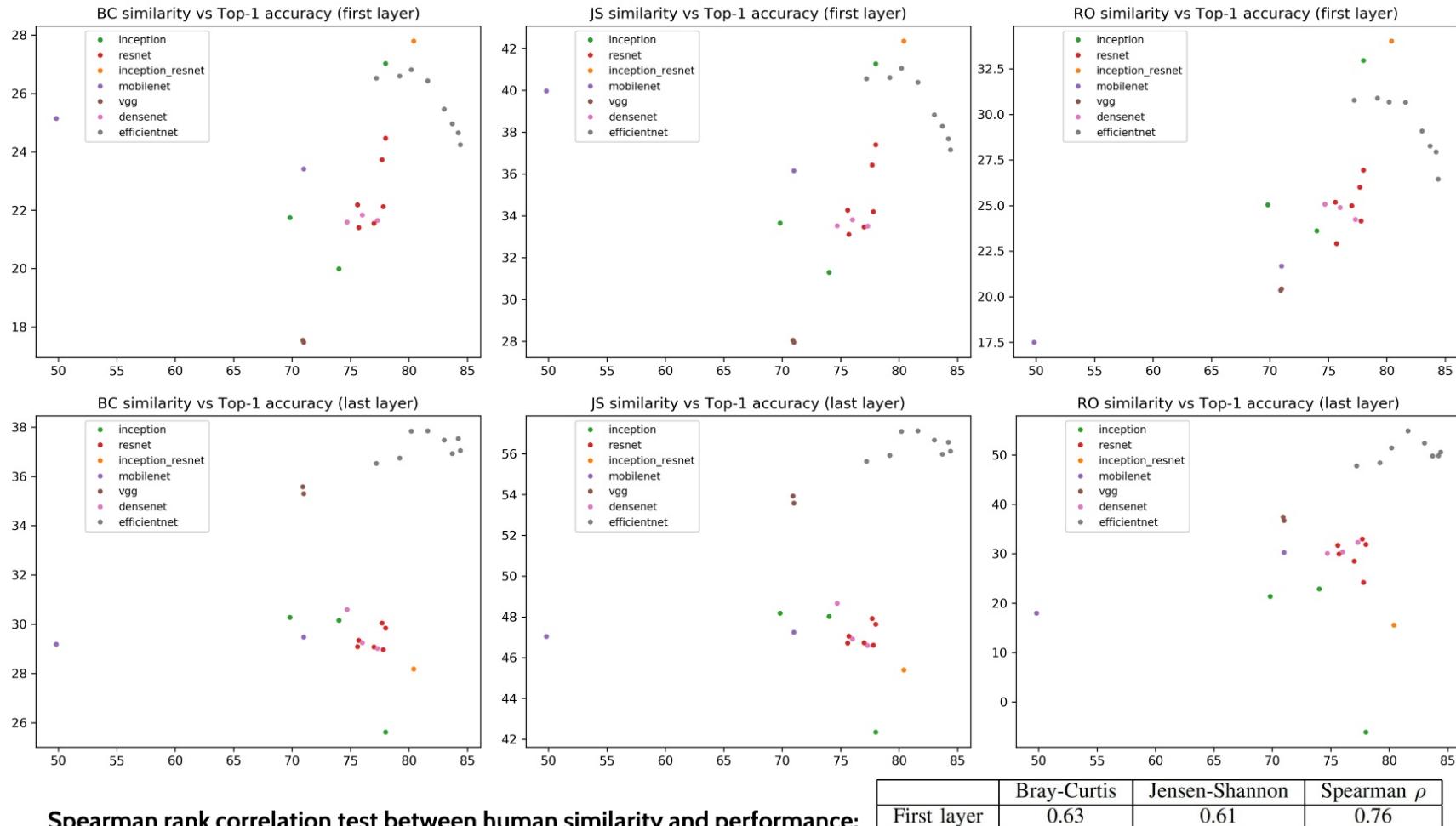




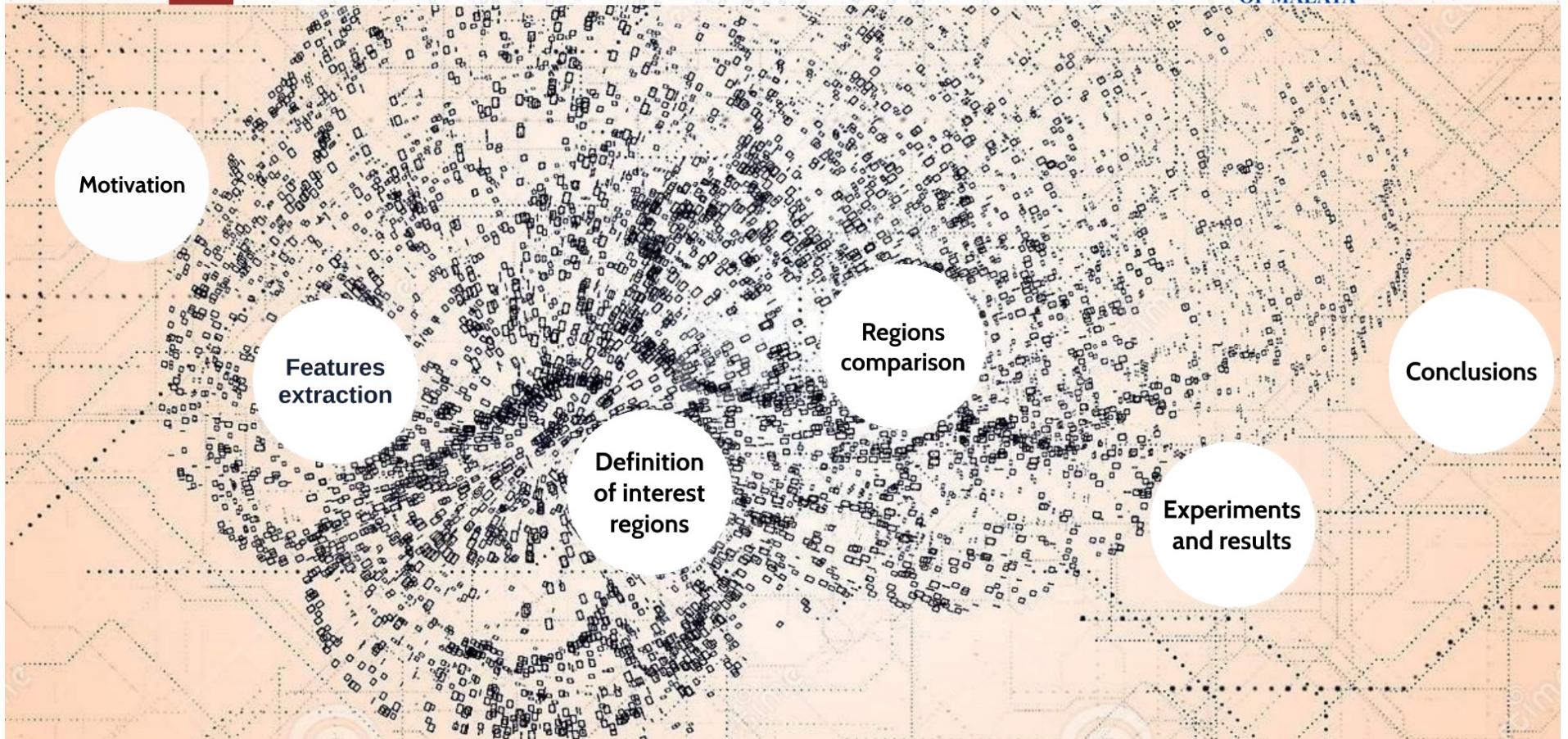
## Summary table of comparisons between human fixations and CNNs

Model	Layer	$BC_{f_{F_j} f_H}$	$JS_{f_{F_j} f_H}$	$\rho_{f_{F_j} f_H}$	$BC_{f_{F_j} f_H}$ R=21.26%	$JS_{f_{F_j} f_H}$ R=34.11%	$\rho_{f_{F_j} f_H}$
		R=21.26%	R=34.11%				
HMAX		15.24%	24.81%	18.47			
AlexNet	$C1$	25.51%	38.95%	31.36	InceptionV3	c1	27.03%
	$C2$	31.61%	48.59%	29.66		c2	26.51%
	$C3$	34.45%	53.05%	37.19		c3	<b>21.87%</b>
	$C4$	34.47%	52.94%	36.38		c4	<b>21.87%</b>
	$C5$	34.09%	52.49%	35.39		c5	<b>21.88%</b>
VGG-19	$C1$	14.03%	22.99%	14.47		c6	27.54%
	$C2$	20.51%	32.12%	21.94		c7	25.62%
	$C3$	20.24%	47.91%	20.34	Densenet-201	$C1$	<b>21.66%</b>
	$C4$	33.46%	50.05%	33.91		$C2$	29.82%
	$C5$	35.18%	53.44%	37.16		$C3$	34.65%
ResnetV2-50	$C1$	22.18%	<b>34.27%</b>	25.20		$C4$	32.36%
	$b1$	32.51%	49.42%	30.24		$C5$	29.02%
	$b2$	35.33%	54.24%	41.12	EfficientNet-b7	$b1$	28.27%
	$b3$	29.42%	47.09%	36.81		$b2$	34.85%
	$b4$	29.09%	46.71%	31.75		$b3$	32.42%
						$b4$	31.21%
						$b5$	33.37%
						$b6$	37.71%
						$b7$	37.37%
							36.69
							40.12
							32.59
							28.63
							31.15
							52.33
							44.92

## Correlation between human similarity and classification performance. x-axis: performance, y-axis: human similarity.



# From early biological models to CNNs: do they look where humans look?



## Conclusions

**Q1.** *Do CNNs look at the same regions humans look at?*

- A. HMAX features and filters responses from shallow CNN's layers are not similar to human fixations, while filter responses of deeper layers are.**

**Q2.** *Is visual attention of recent networks more similar to humans attention than the attention of earlier developed networks?*

- A. It is if we consider early models such as HMAX, while the attention of networks such as AlexNet and VGG is more or less as similar to humans as that of more recent networks such as Densenet and Efficientnet.**

**Q3.** *Is CNNs-human visual attention similarity related to CNNs classification performance?*

- A. Yes it is: there is positive correlation between similarity and performance.**

The image features a large, bold red word "thank you" at the center. Surrounding it are numerous other words in various languages, each with its own unique color and size, creating a diverse and colorful collage. The surrounding words include: "спасибо" (spasibo) in red, "danke" in blue, "谢謝" (Xie Xie) in light blue, "ngiyabonga" in orange, "teşekkür ederim" in pink, "gracias" in green, "mochchakkeram" in blue, "go raibh maith agat" in purple, "dakujem" in pink, "merci" in red, "감사합니다" (Gamsahamnida) in yellow, "terima kasih" in yellow, "arigatō" in pink, "dank je" in pink, "obrigado" in green, "dziekuje" in pink, "sukriya" in purple, "khop khun krap" in green, "grazie" in blue, "merci" in red, and "спасибо" (spasibo) again in red.

# From early biological models to CNNs: do they look where humans look?

