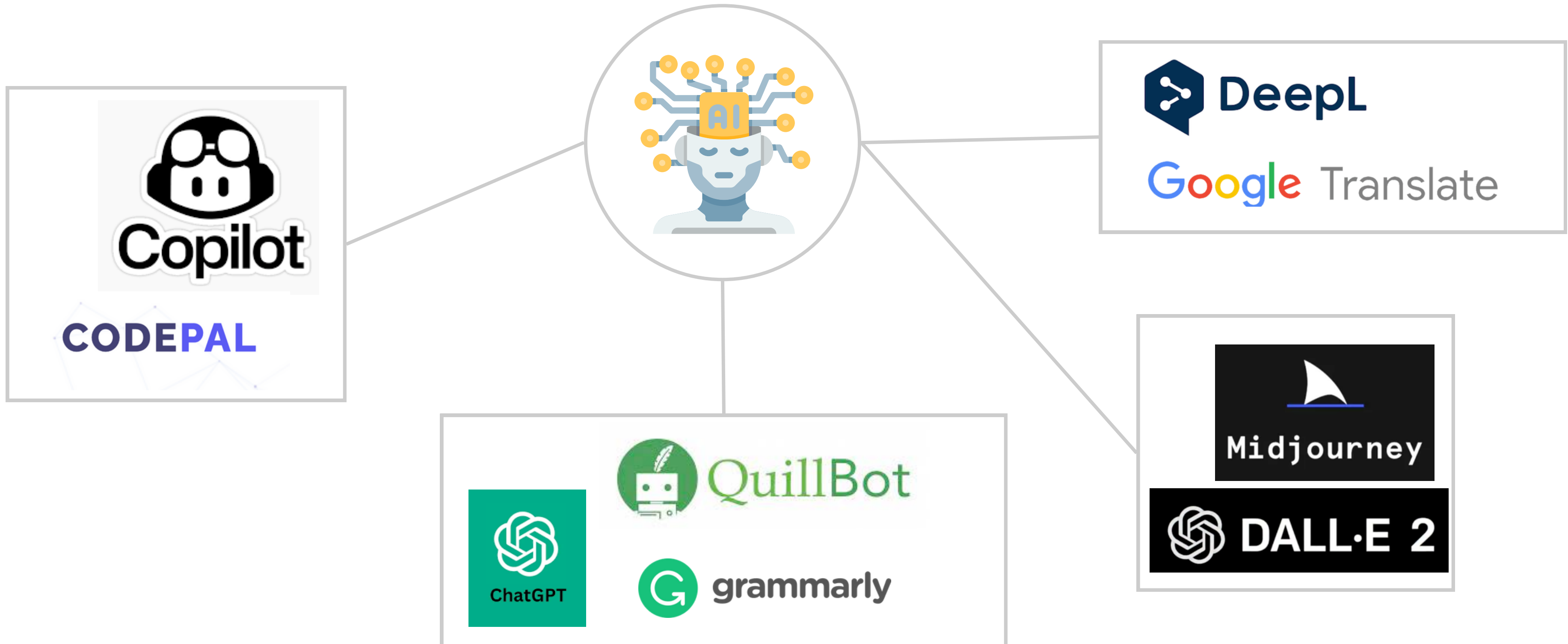


A dark, grayscale background image showing a person's hands holding a smartphone. The hands are positioned as if they are about to interact with the device. The image is slightly blurred and has a dark overlay, making it a subtle background for the text.

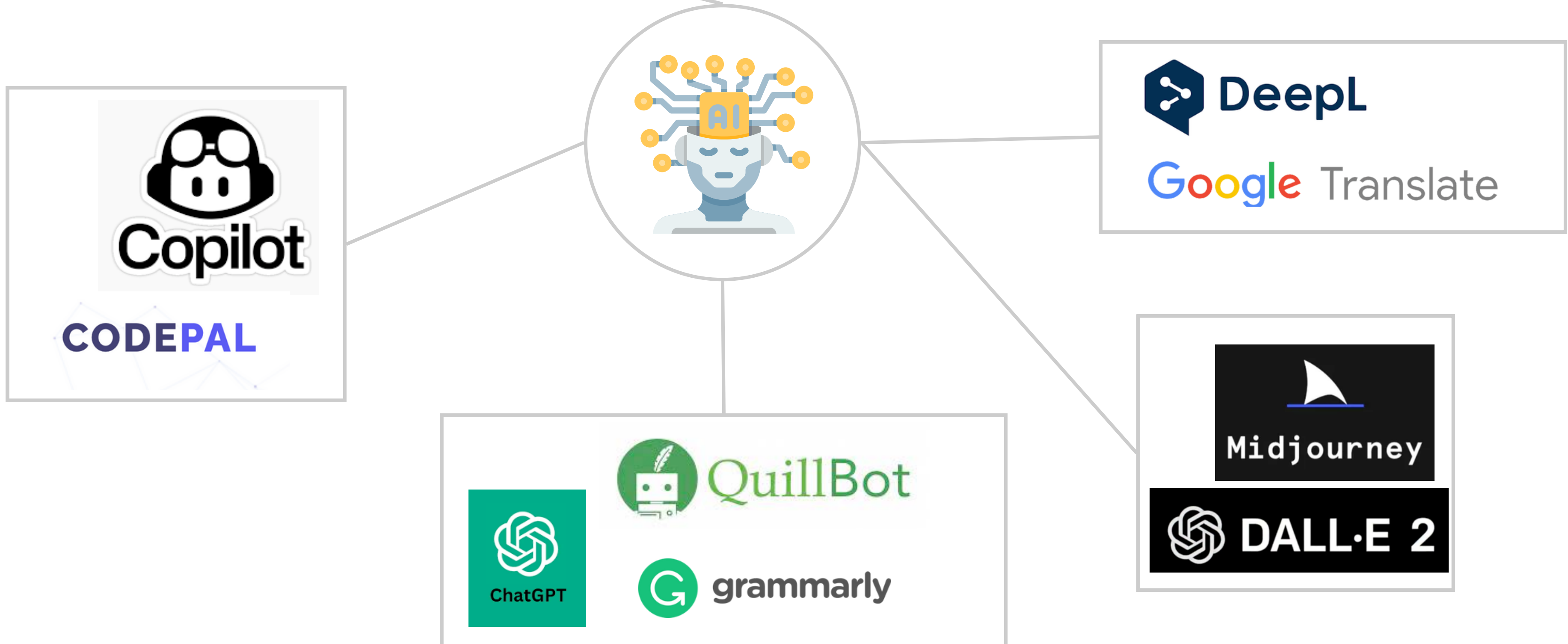
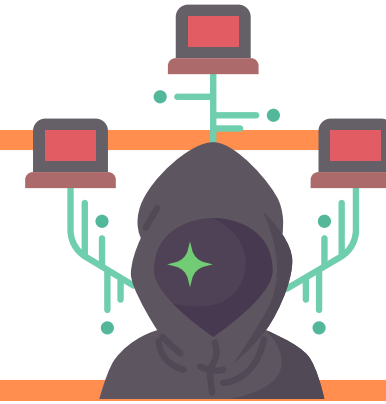
Breaking Down AI to get Explanations

Karel Hynek

AI is everywhere



Computer Security





Me!

- **PhD Candidate** at FIT Czech Technical University
- Four years developing **Network security detectors**
- Mainly focusing on **AI-based security detectors**

Accuracy of

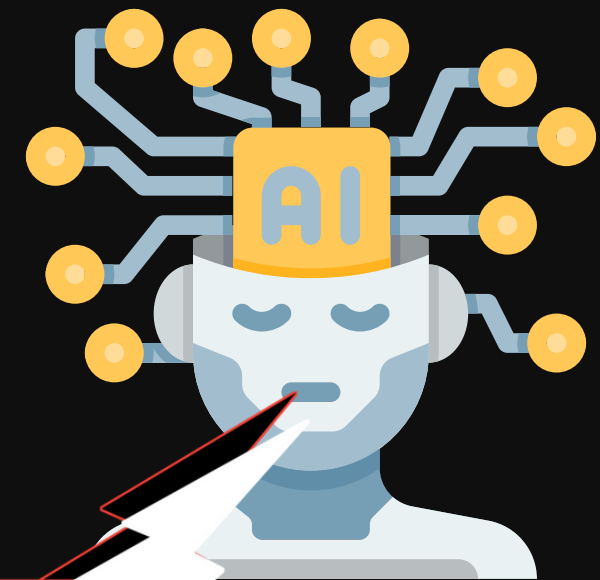
99.9%

**But it is hard to verify the
output.**

Are those hotdogs?



Are those hotdogs?



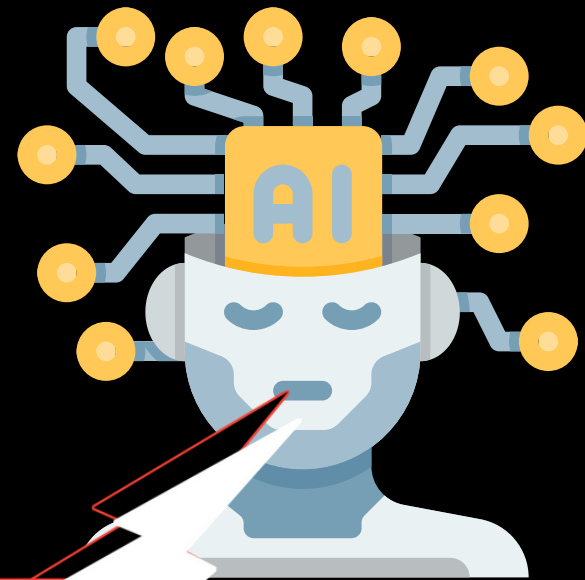
YES!

Is this malware?

Packet Sizes: [231|412|221|82|232|415|240|82|242|414|230|82|237|415|244|82|261|75|412|64],
Interpacket Times: [10 | 14 | 443|13|14 | 53 | 534|44|424|232|535|44|232| 1 | 13 | 44|434|32|555|23]
Directions: [1 | -1 | 1 | 1 | 1 | 1 | -1 | -1 | 1 | -1 | 1 | 1 | -1 | -1 | -1 | -1 | 1 | -1 | 1 | 1]

Is this malware?

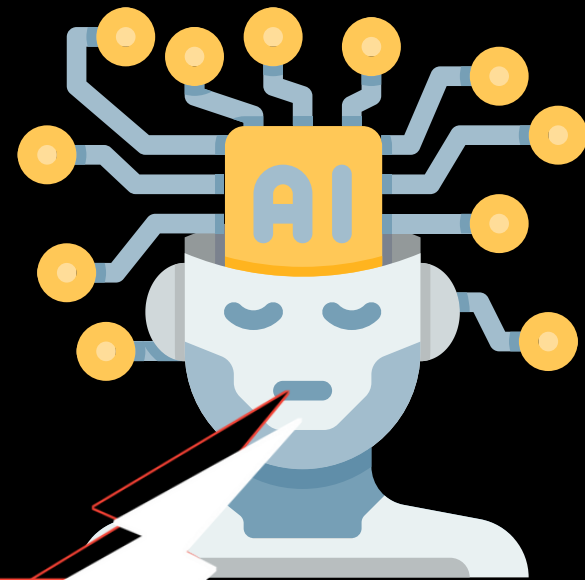
Packet Sizes: [231|412|221|82|232|415|240|82|242|414|230|82|237|415|244|82|261|75|412|64]
Interpacket Times: [10 | 14 | 443|13|14 | 53 | 534|44|424|232|535|44|232| 1 | 13 | 44|434|32|555|23]
Directions: [1 | -1 | 1 | 1 | 1 | 1 | -1 | -1 | 1 | -1 | 1 | 1 | -1 | -1 | -1 | -1 | 1 | -1 | 1 | 1]



YES!

Is this malware?

Packet Sizes: [231|412|221|82|232|415|240|82|242|414|230|82|237|415|244|82|261|75|412|64]
Interpacket Times: [10 | 14 | 443|13|14 | 53 | 534|44|424|232|535|44|232| 1 | 13 | 44|434|32|555|23]
Directions: [1 | -1 | 1 | 1 | 1 | 1 | -1 | -1 | 1 | -1 | 1 | 1 | -1 | -1 | -1 | -1 | 1 | -1 | 1 | 1]



YES!

**You can't quickly
verify.**

We must trust it...

99.9%

Quick Math...

10K Network
telemetry
records per
second

99.9%
Accuracy

1 false-
detection
per second

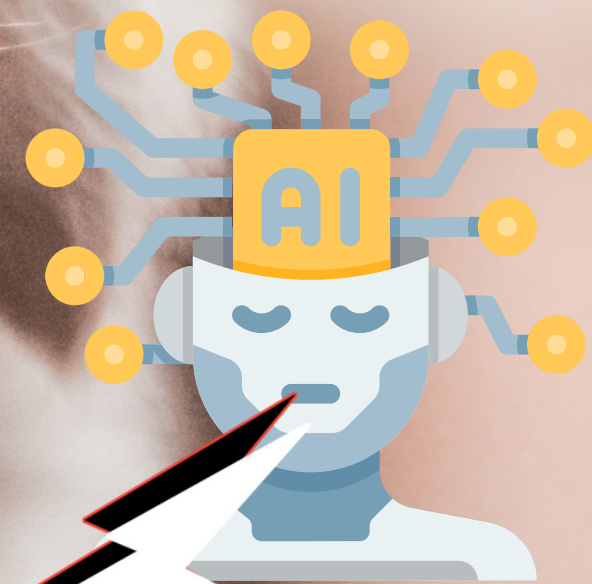
Network of mid-
sized company

Higher standard for
network threat
classification

$10K * 0.001$

It is gonna be turned off very
soon!



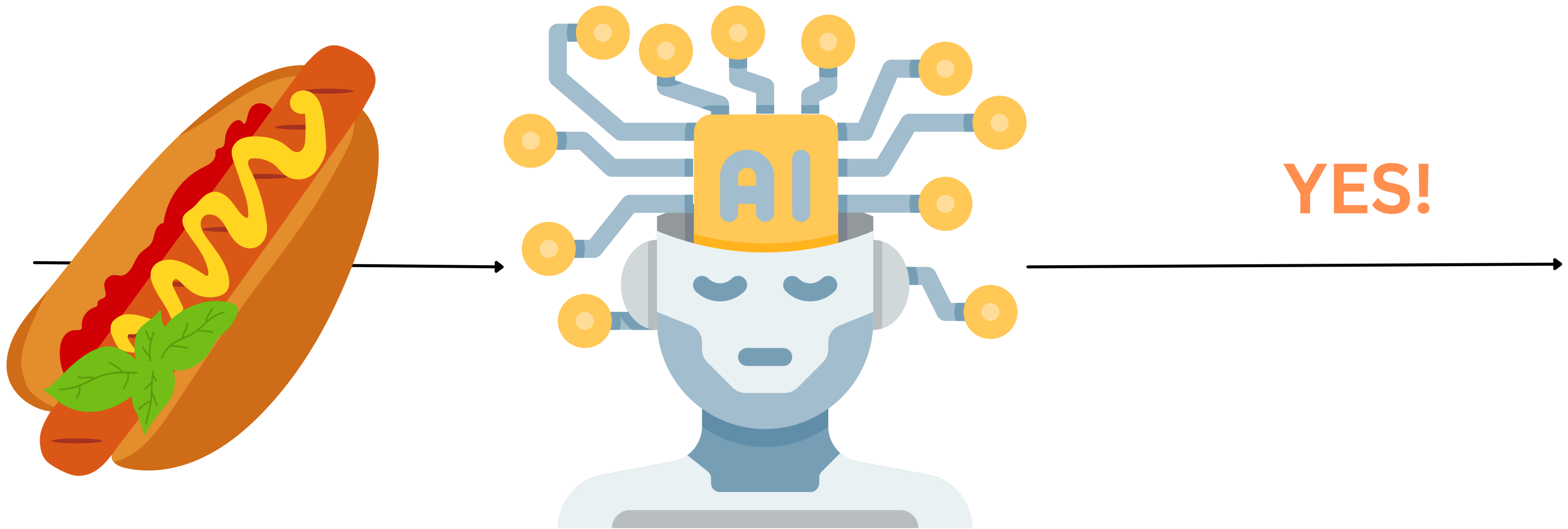


Hotdog!

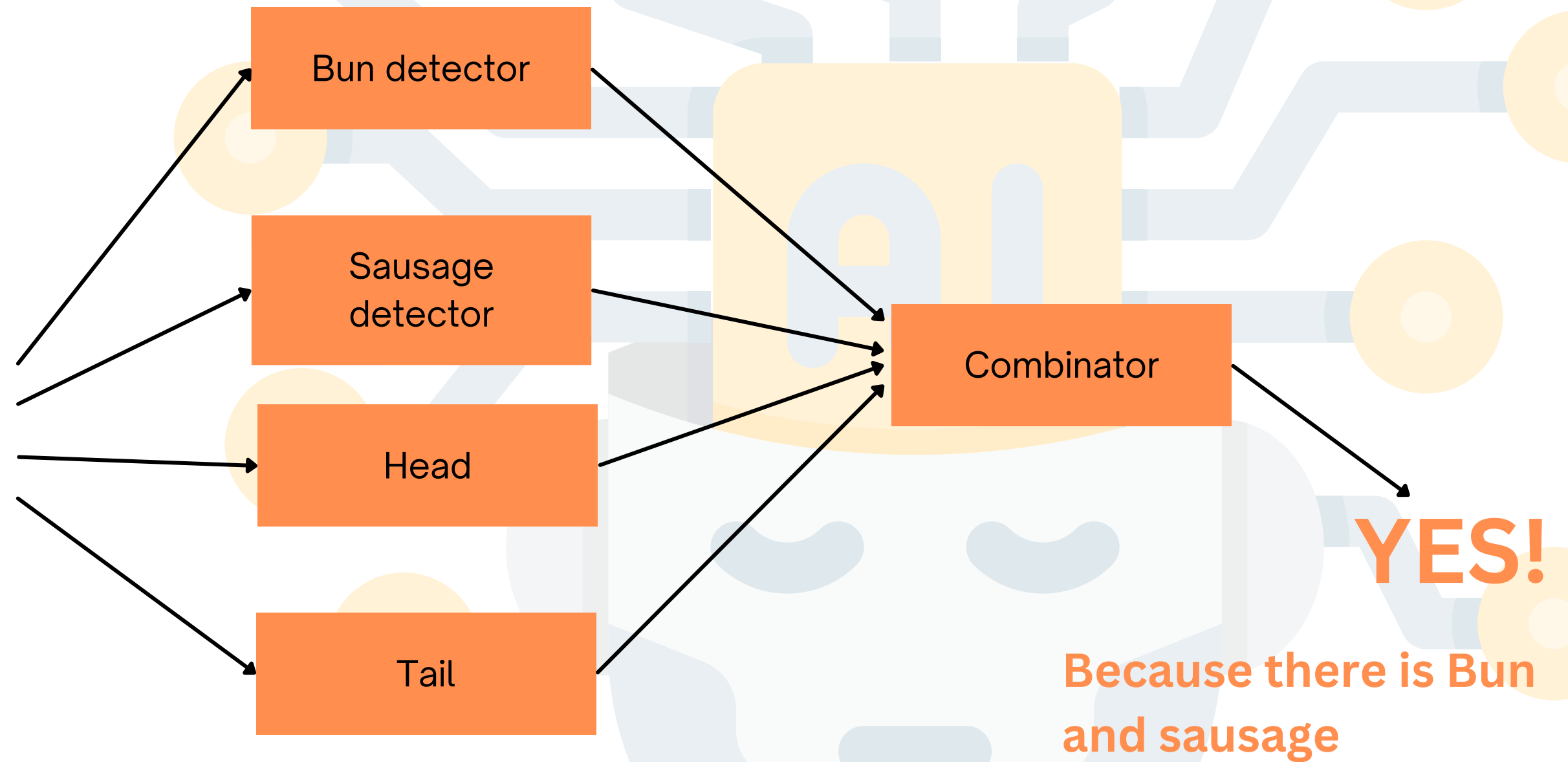
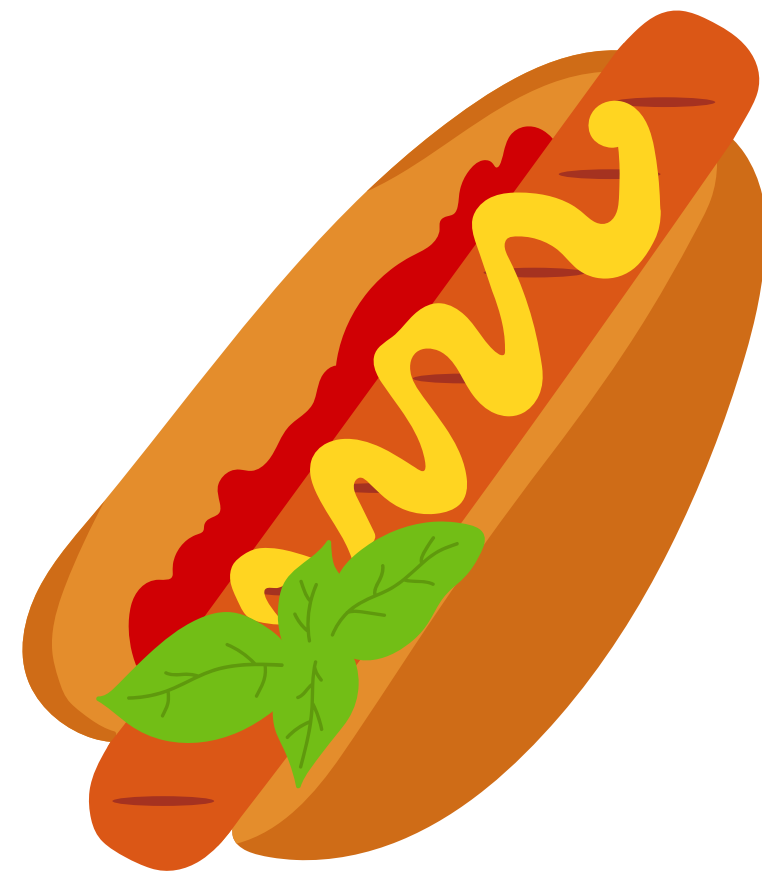
**We need to
distinguish wrong
detections
quickly...**

By adding explanations.

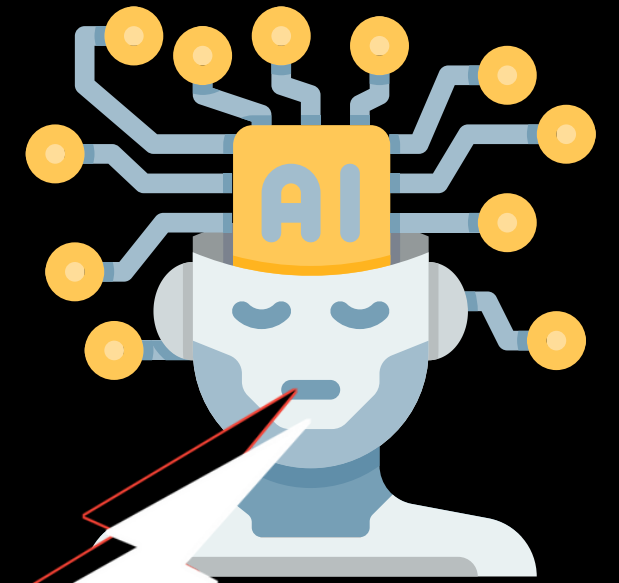
Monolithic AI



Dividing AI



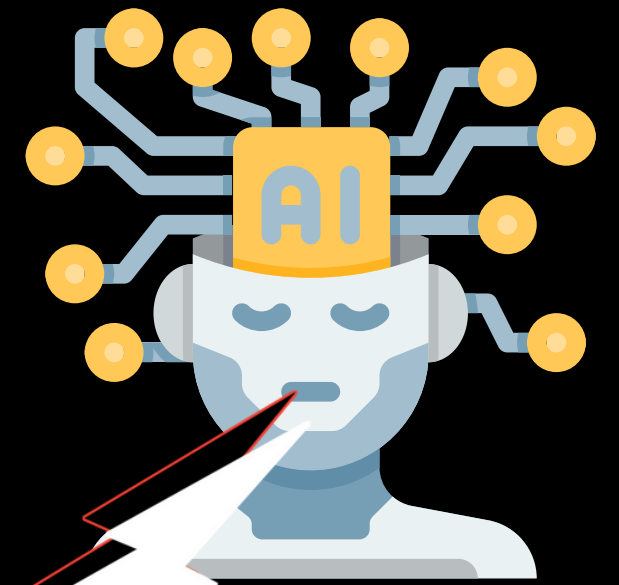
Is this a hotdog?



YES!

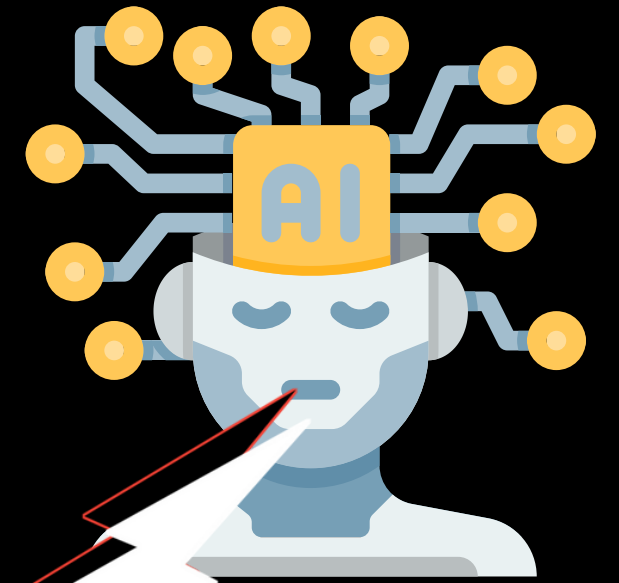
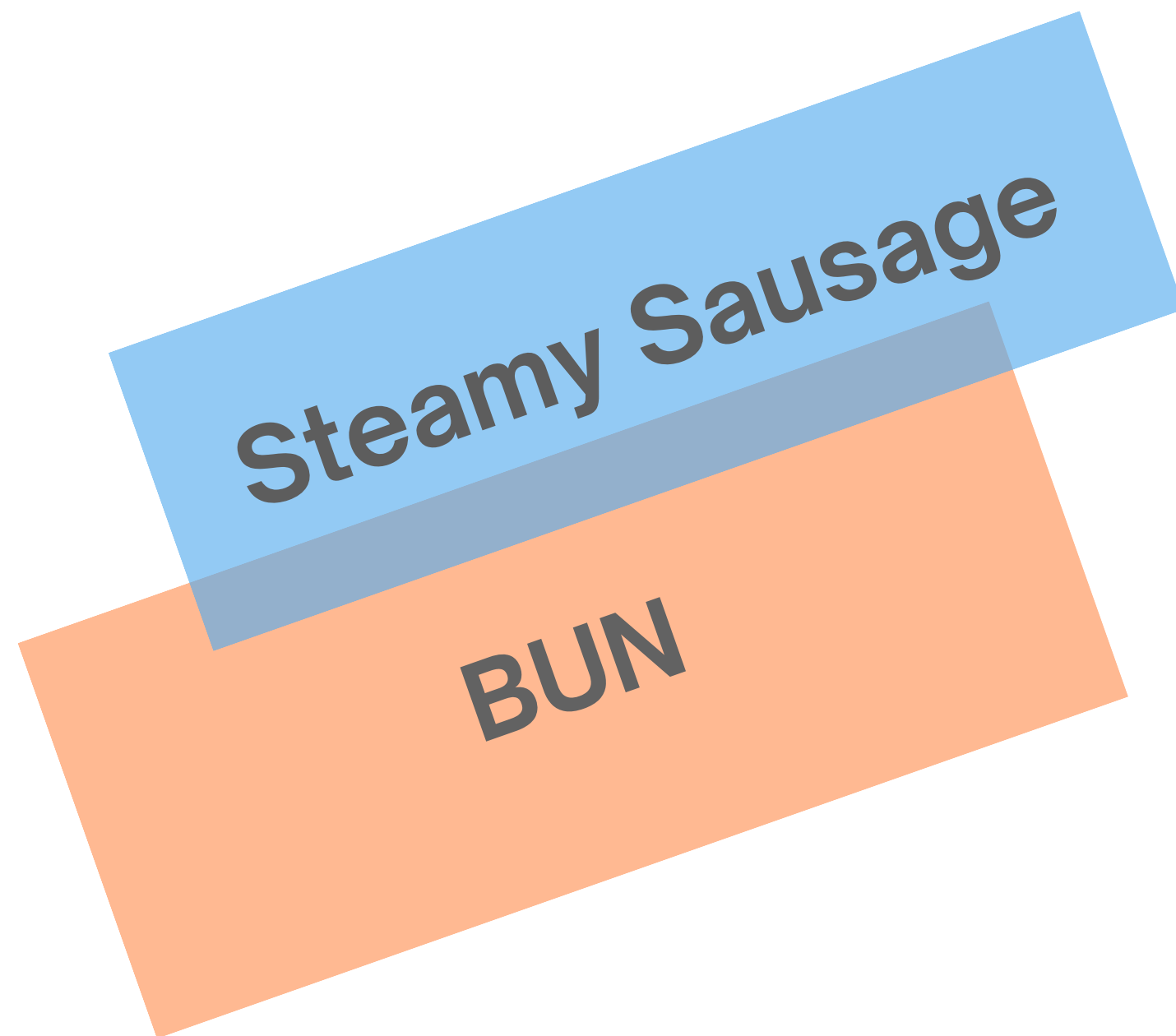
Is this a hotdog?

BUN



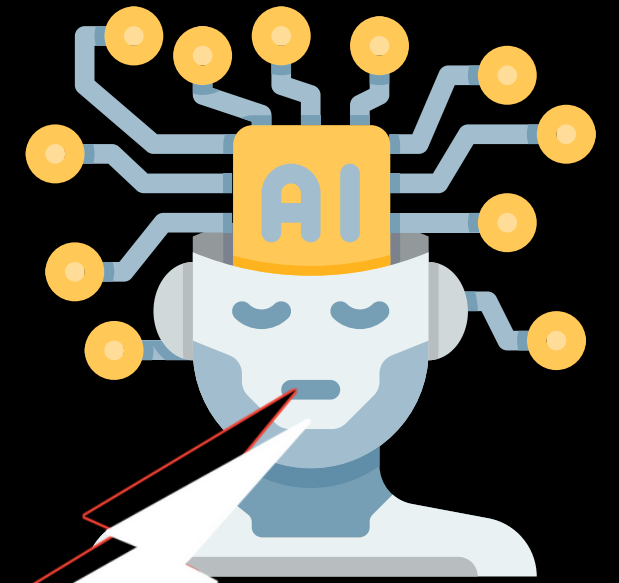
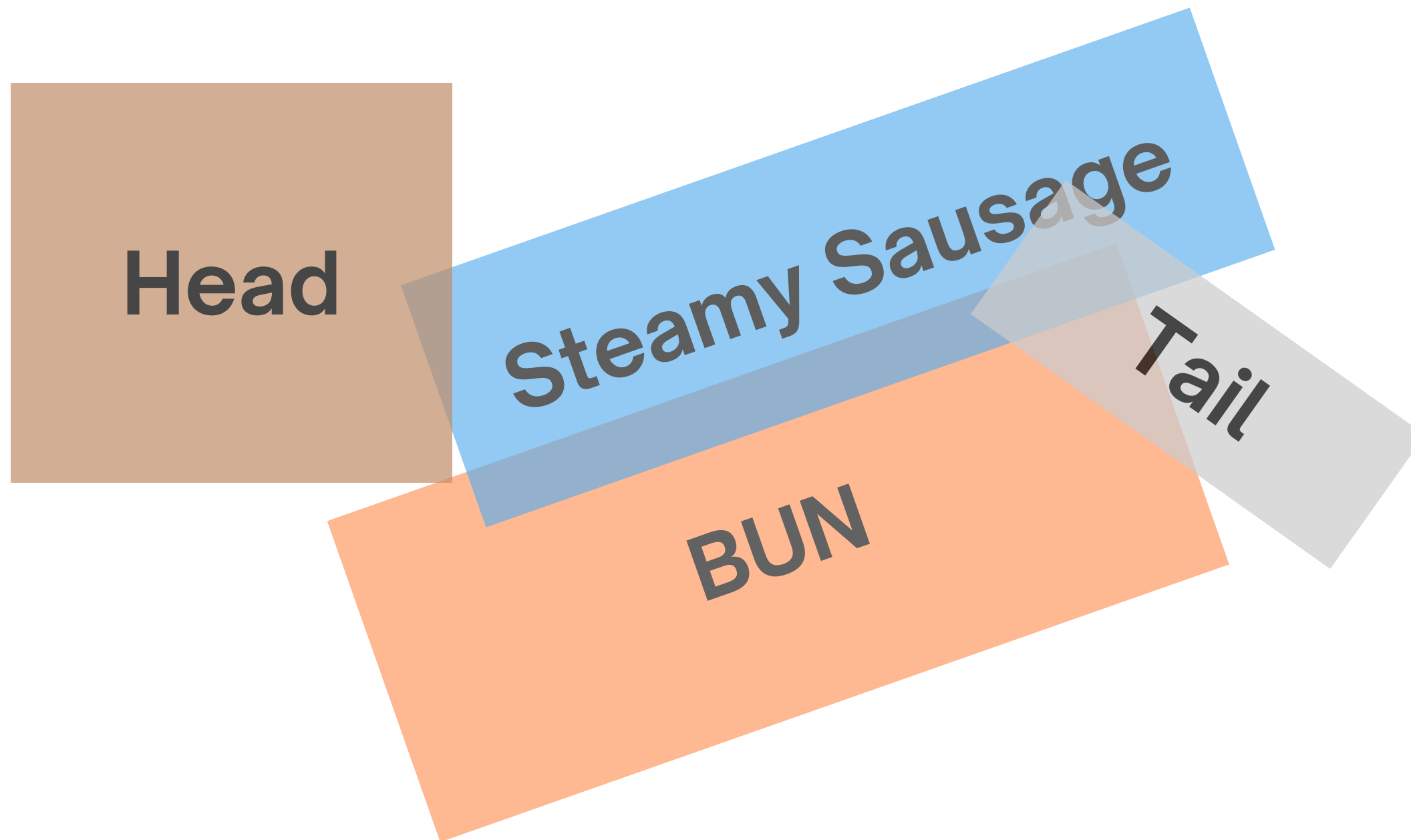
YES!

Is this a hotdog?



YES!

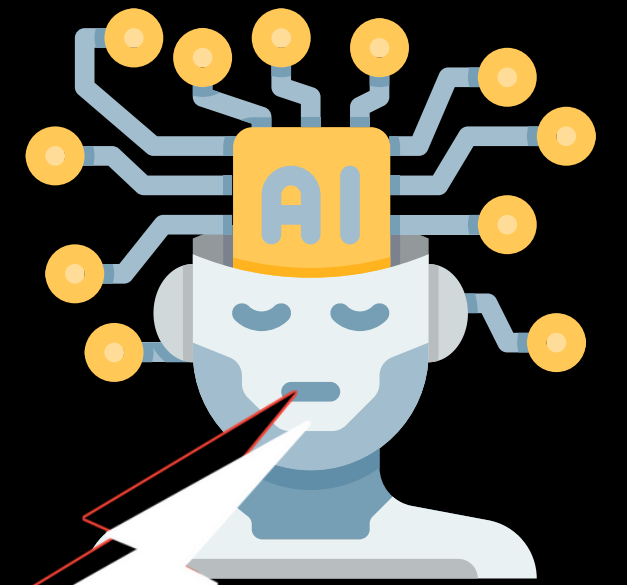
Is this a hotdog?



YES!

**Something is
really wrong!**

Is this a hotdog?



YES!

**Something is
really wrong!**

Dividing AI

[231|412|221|82|232|415|240|82|242|414|230|8
2|237|415|244|82|261|75|412|64]
[10 | 14 |443|13|14 | 53
|534|44|424|232|535|44|232| 1 | 13
|44|434|32|555|23]
[1 | -1 | 1 | 1 | 1 | 1 | -1 | -1 | 1 | -1 | 1 | 1 | -1 | -1 | -1
| -1 | 1 | -1 | 1 | 1]

Periodic
Communication

Covert
Communication

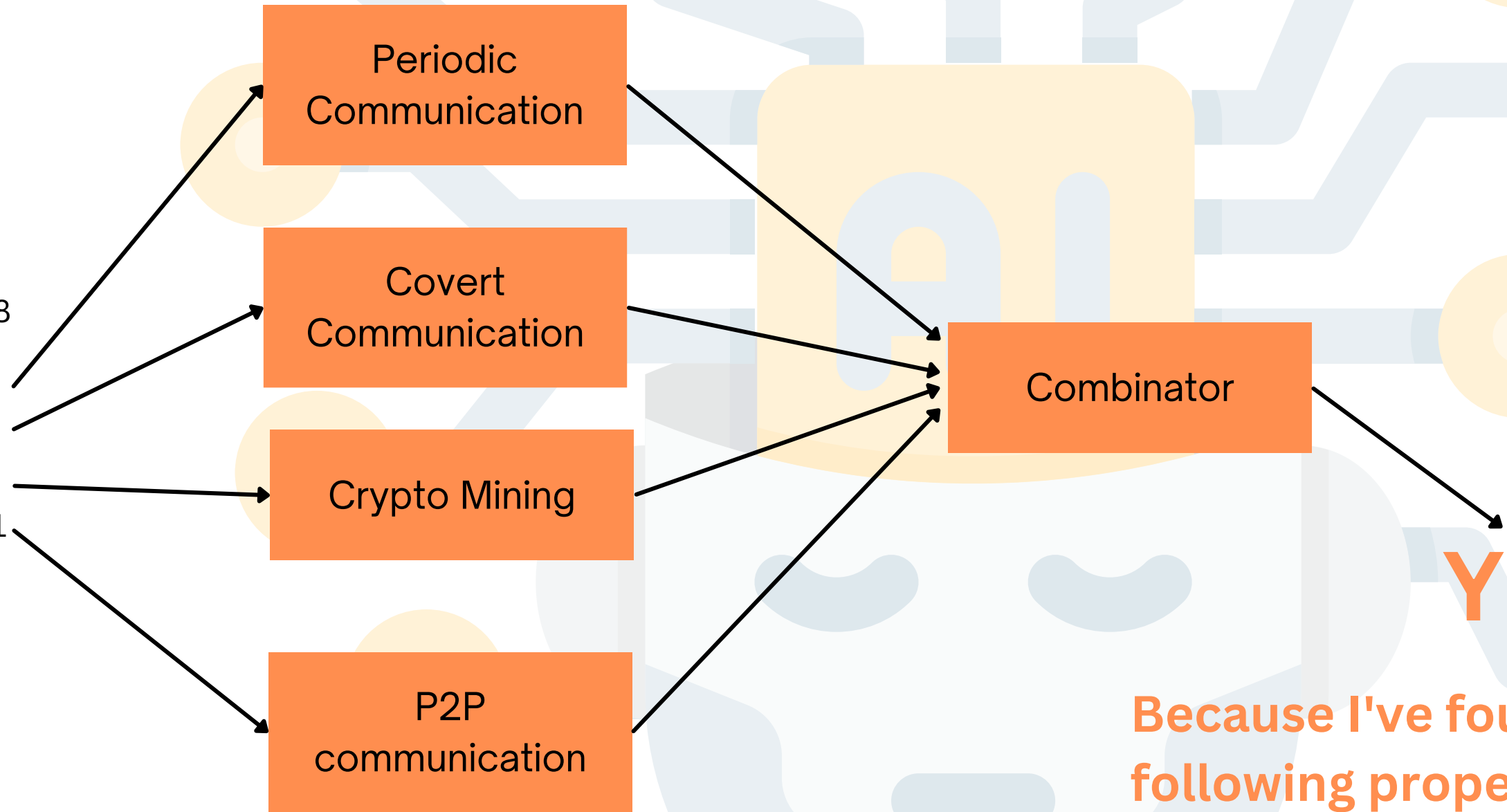
Crypto Mining

P2P
communication

Combinator

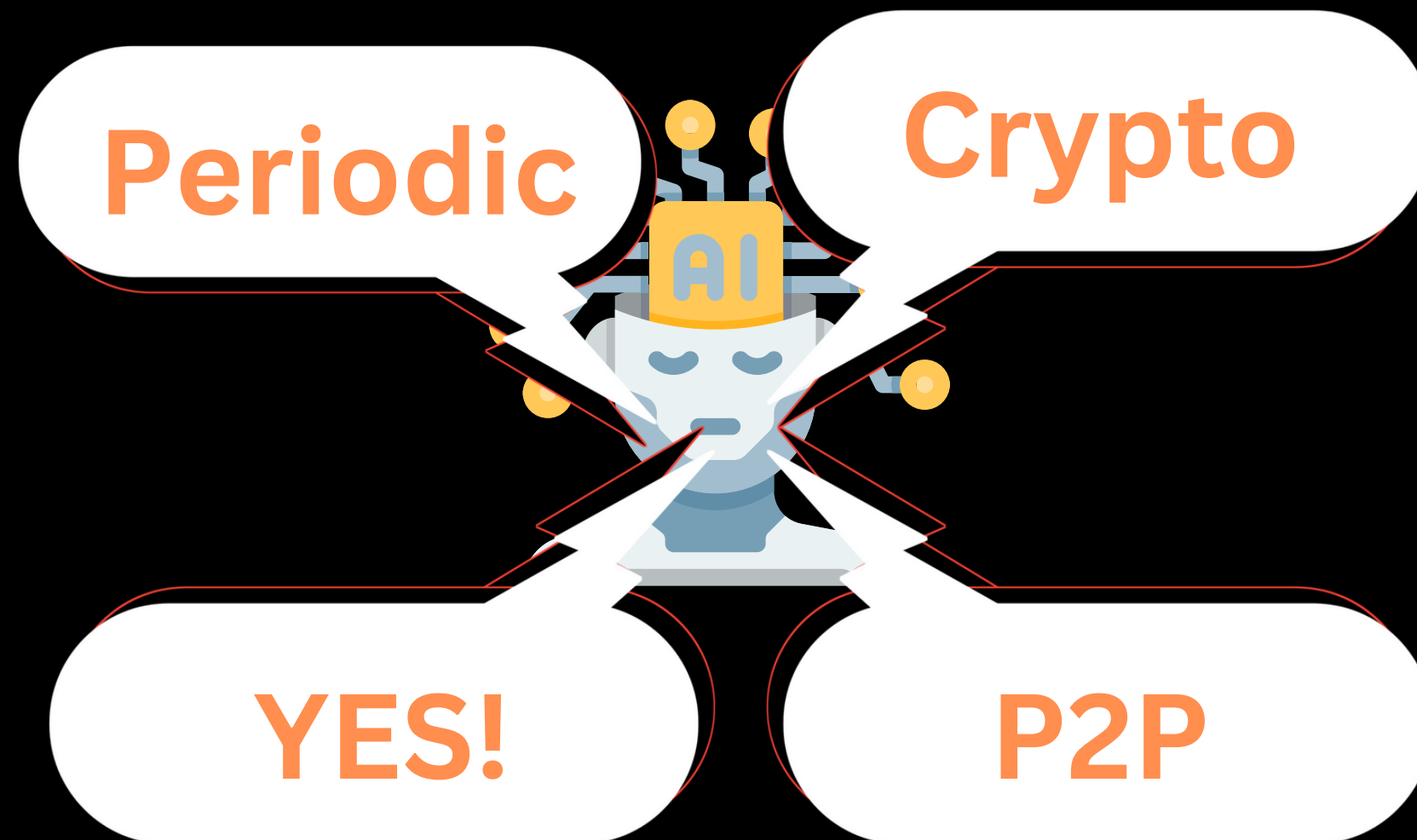
YES!

Because I've found
following properties...



Is this communication malicious?

Packet Sizes: [231|412|221|82|232|415|240|82|242|414|230|82|237|415|244|82|261|75|412|64]
Interpacket Times: [10 | 14 | 443|13|14 | 53 | 534|44|424|232|535|44|232| 1 | 13 | 44|434|32|555|23]
Directions: [1 | -1 | 1 | 1 | 1 | 1 | -1 | -1 | 1 | -1 | 1 | 1 | -1 | -1 | -1 | -1 | 1 | -1 | 1 | 1]



The **benefits** of divided AI

- 1 We can quickly identify misclassifications**
- 2 Less prone to design errors**
- 3 Even non-AI expert can understand the explanations**

**You can
contact me at**

karel.hynek@cesnet.cz

Karel Hynek

Network Security Research and Development
CESNET a.l.e.