



UiT The Arctic University of Norway

NORA summer school on multimodal learning

# Fundamentals of multimodal classification

Rwiddhi Chakraborty

*UiT Machine Learning Group and Visual Intelligence*

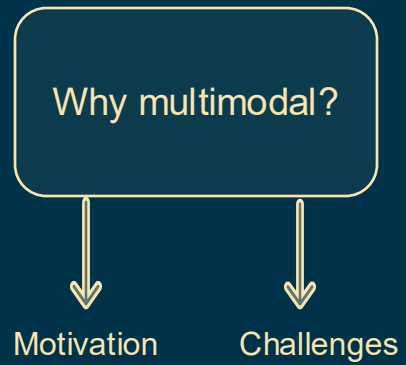
# In this talk

Why multimodal?

Multimodal Chef

Modern practices

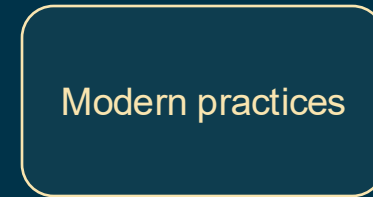
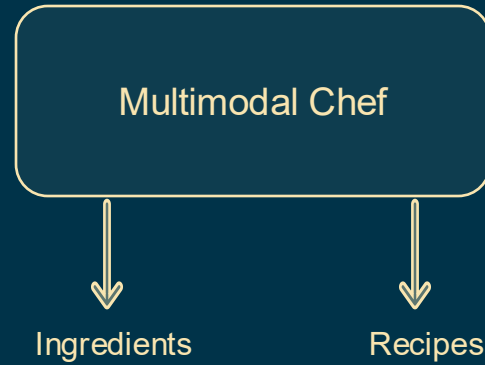
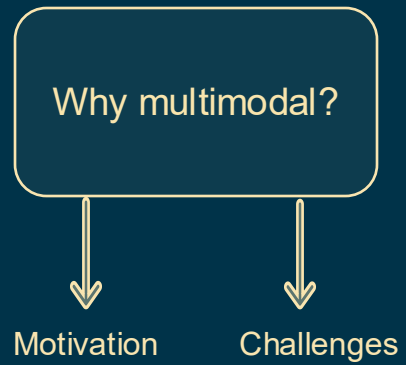
# In this talk



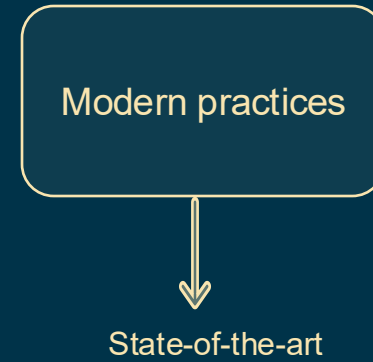
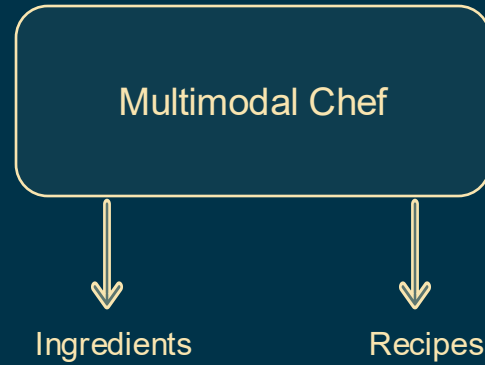
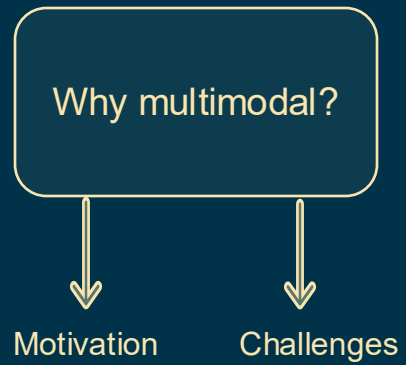
Multimodal Chef

Modern practices

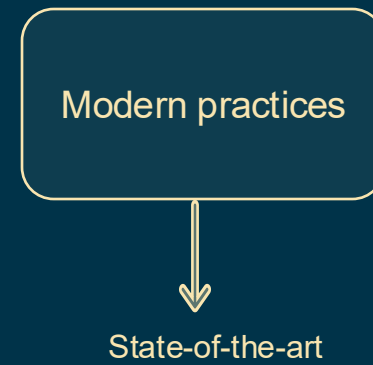
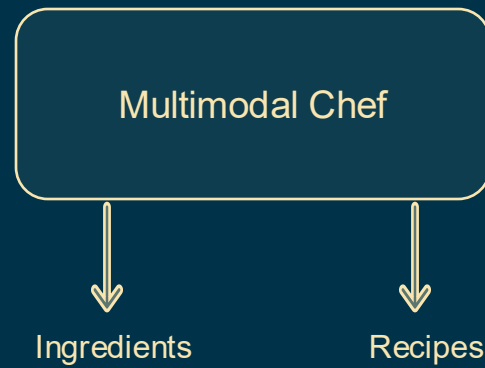
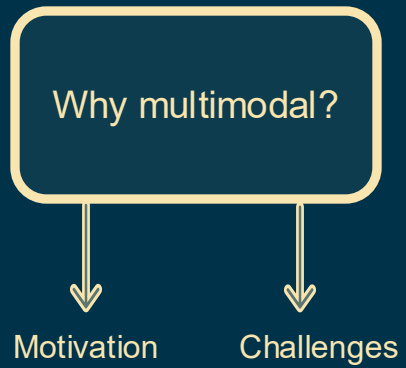
# In this talk



# In this talk



# In this talk



# Why Multimodal?

Data exists in many forms

# Why Multimodal?





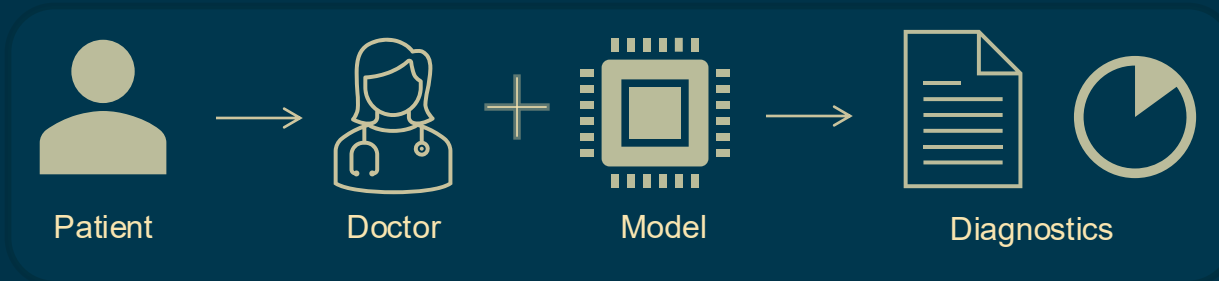
# Why Multimodal?



Diversity of forms  $\Rightarrow$  diversity of information  $\Rightarrow$  better learning

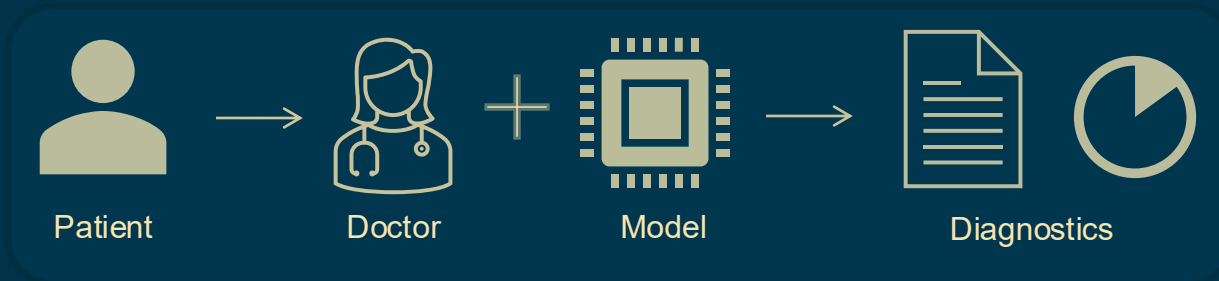
# In the wild

*Patient X visits an AI-assisted health service to generate a diagnostic report.*



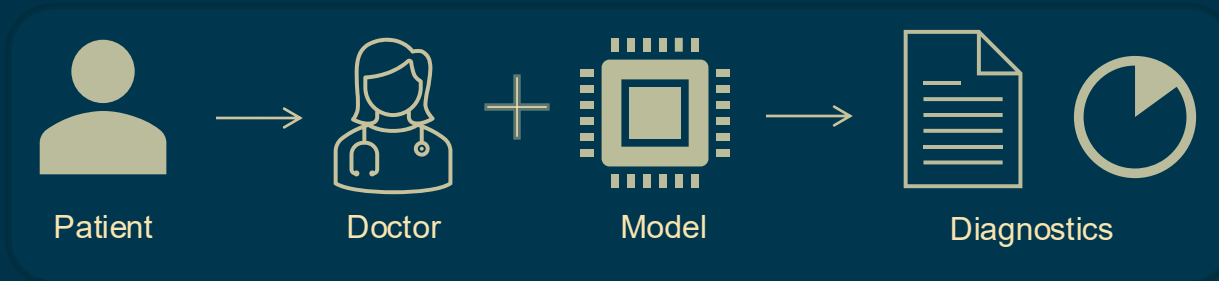
# In the wild

*Patient X visits an AI-assisted health service to generate a diagnostic report.  
What kind of data can we access?*



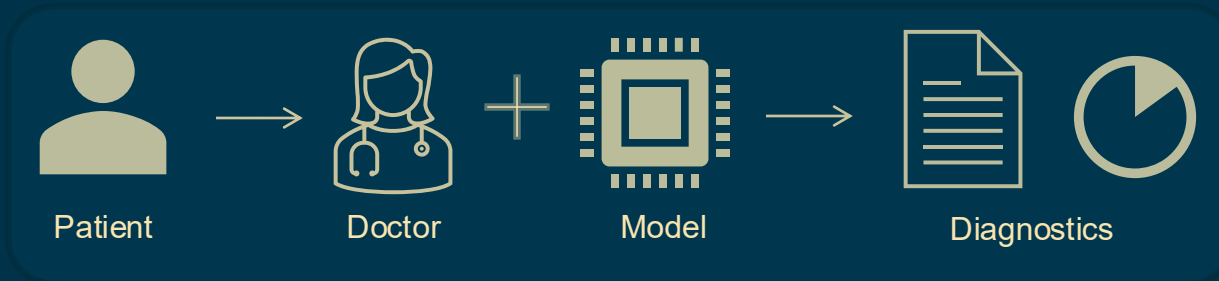
# In the wild

*Patient X visits an AI-assisted health service to generate a diagnostic report.  
What kind of data can we access?*



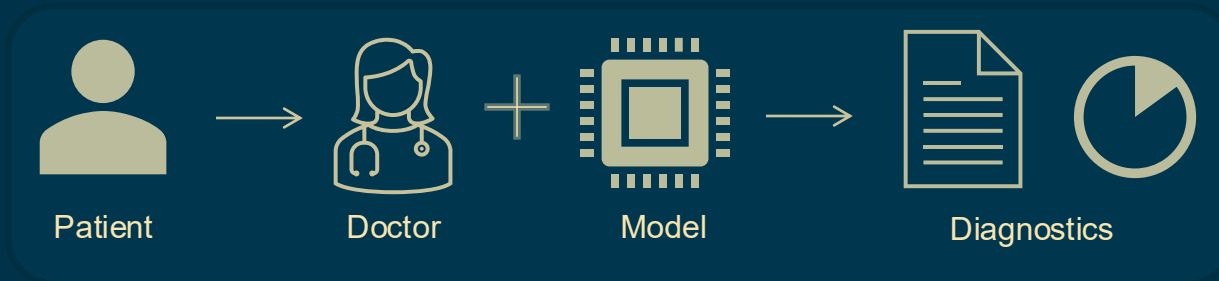
# In the wild

*Patient X visits an AI-assisted health service to generate a diagnostic report.  
What kind of data can we access?*

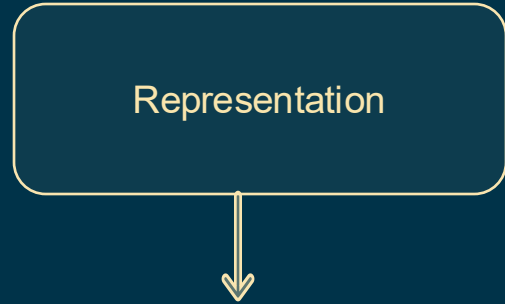


# In the wild

*Patient X visits an AI-assisted health service to generate a diagnostic report.  
What kind of data can we access?*



# Challenges



How do we extract relevant features?

# Challenges

Representation



How do we extract relevant features?

Combination



How do we fuse representations?



# Challenges

Representation



How do we extract relevant features?

Combination



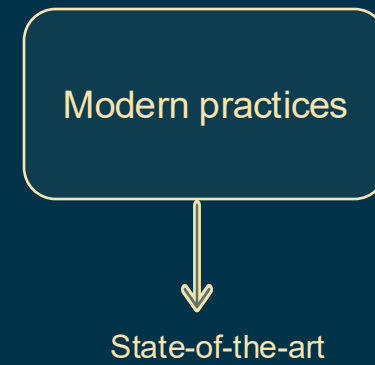
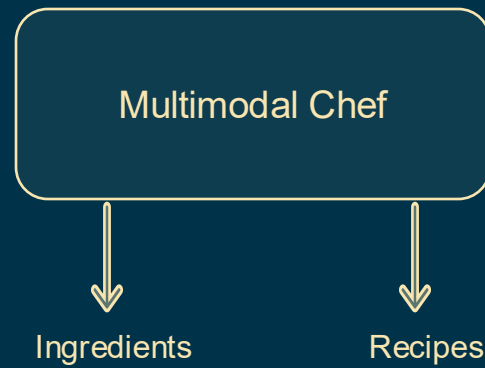
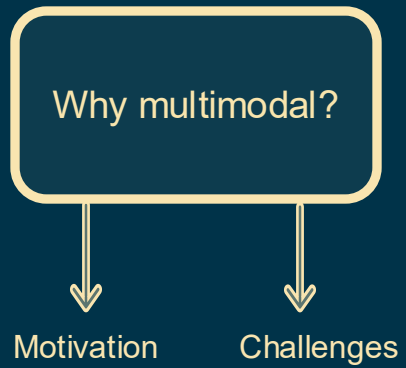
How do we fuse representations?

Labels



How do we learn with limited supervision?

# In this talk



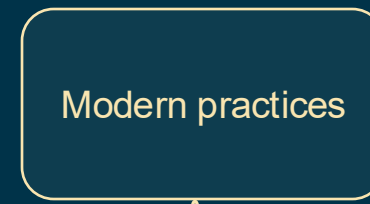
# In this talk



Motivation      Challenges

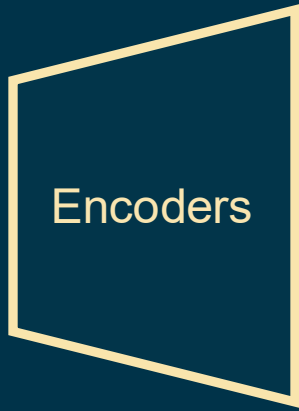


Ingredients      Recipes



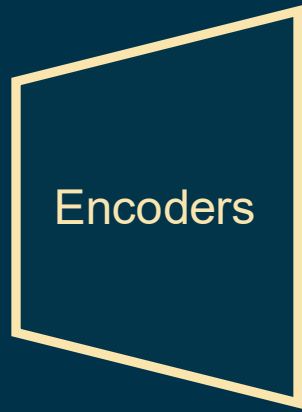
State-of-the-art

# Multimodal Chef: The Ingredients

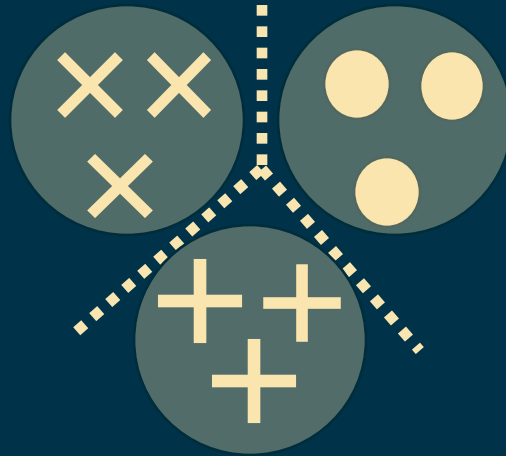


Feature Extraction

# Multimodal Chef: The Ingredients

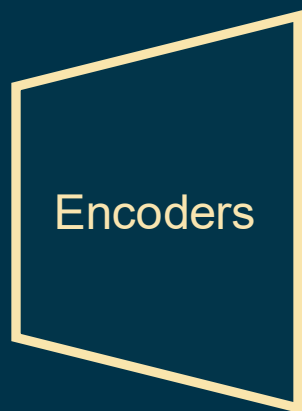


Feature Extraction

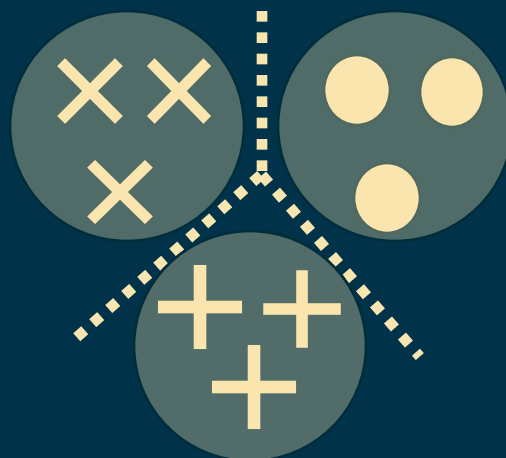


Representations

# Multimodal Chef: The Ingredients



Feature Extraction

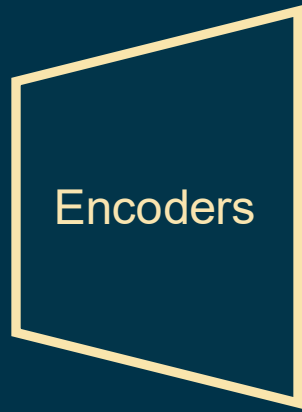


Representations

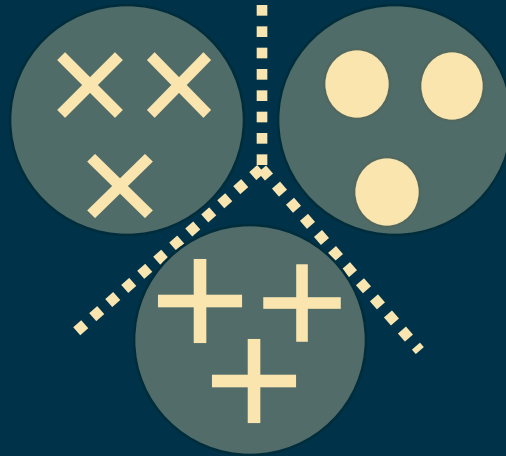


Fuse & Learn

# Multimodal Chef: The Ingredients



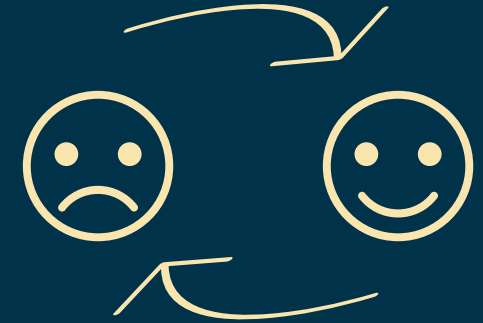
Feature Extraction



Representations

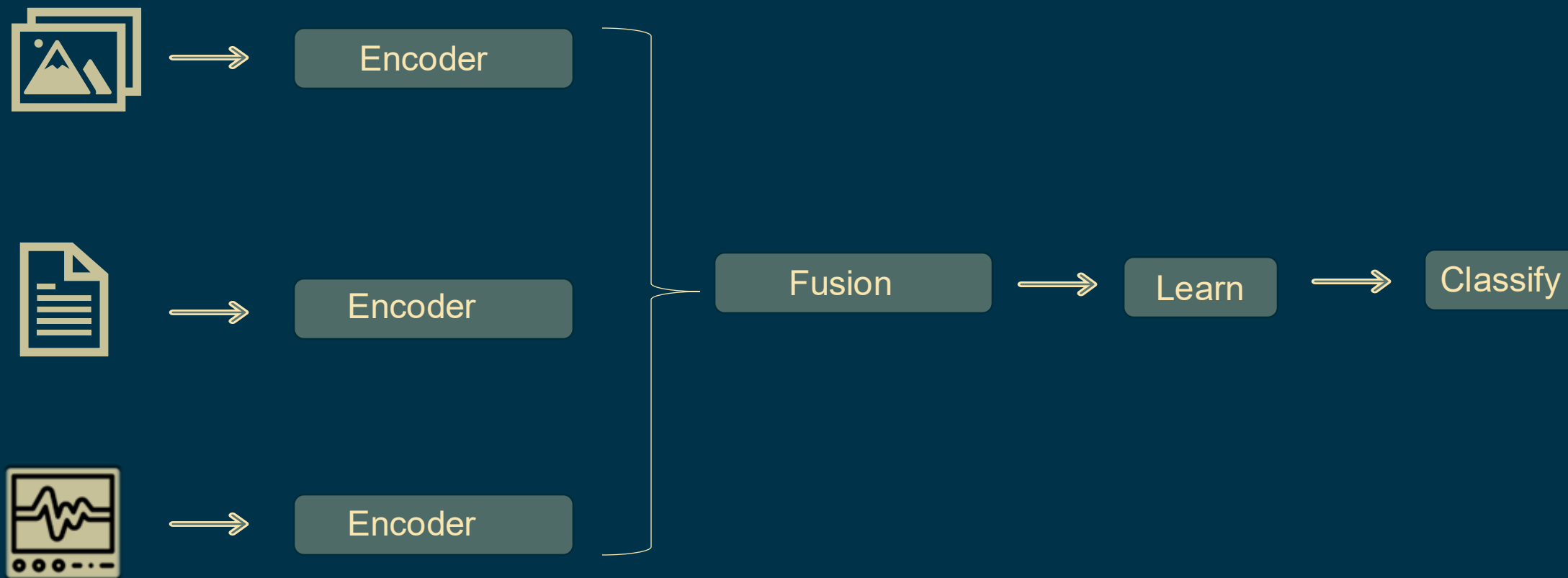


Fuse & Learn



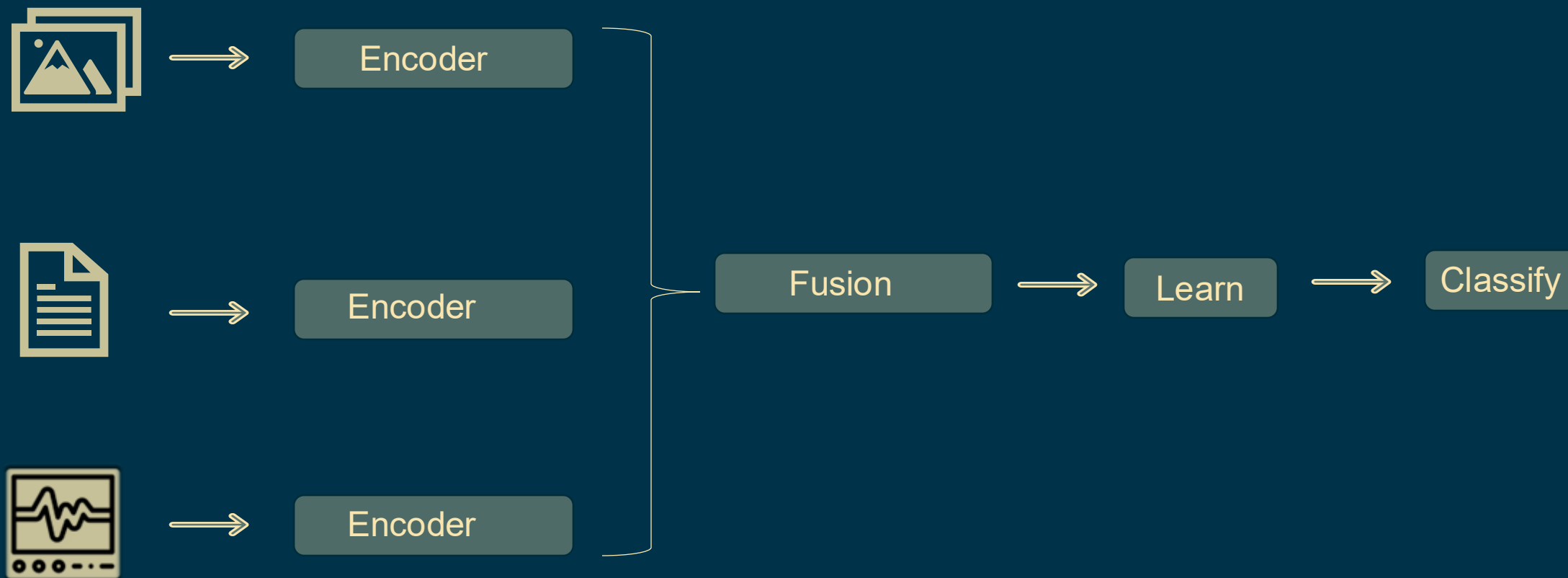
Classify

# Multimodal Chef: The Ingredients

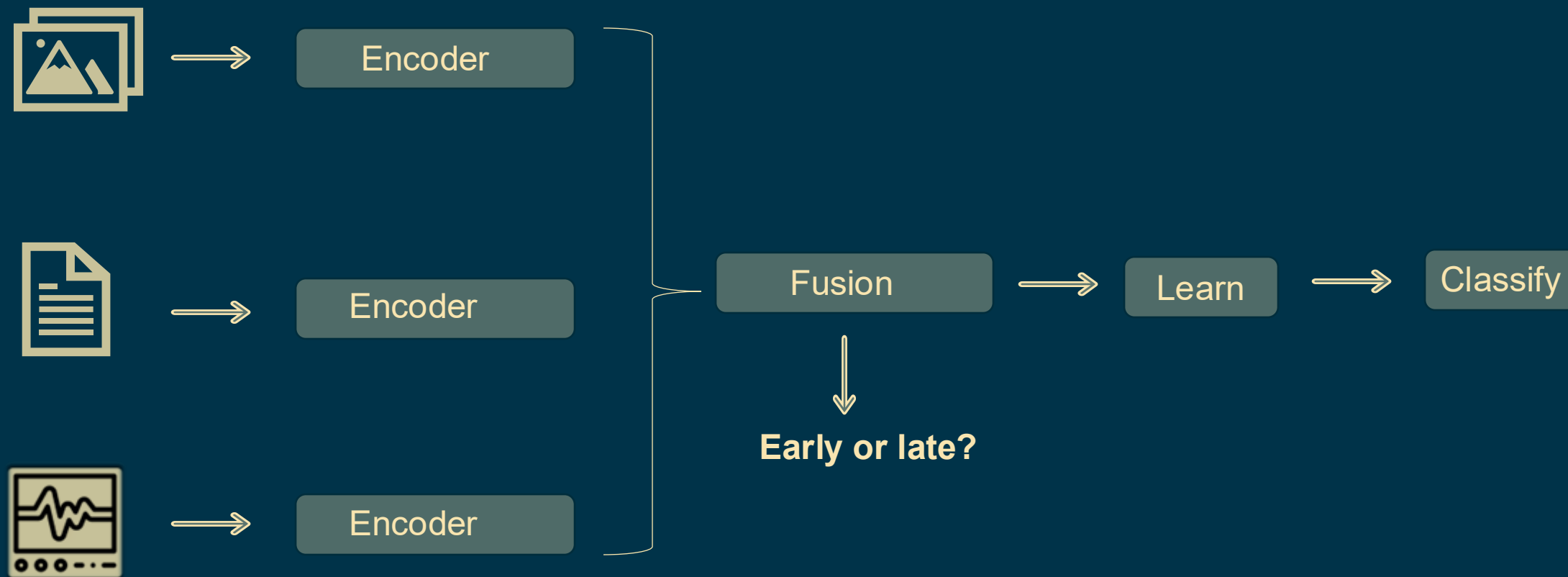




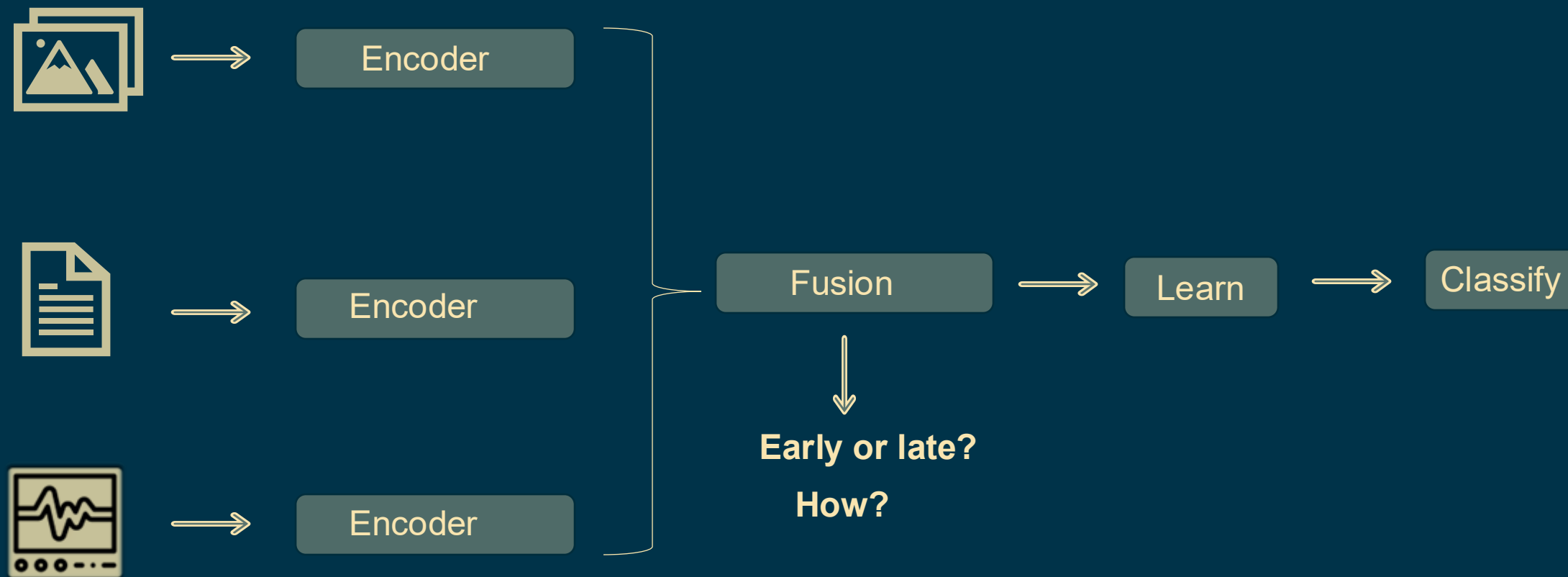
# Multimodal Chef: Recipes



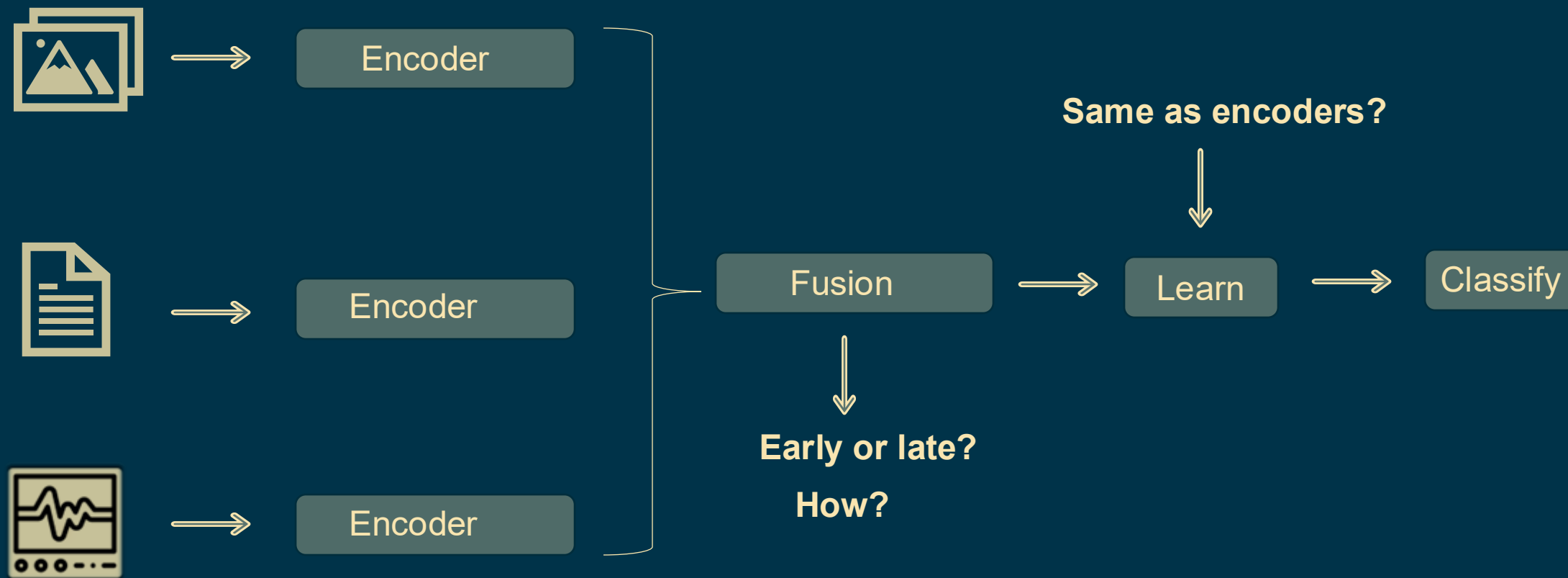
# Multimodal Chef: Recipes



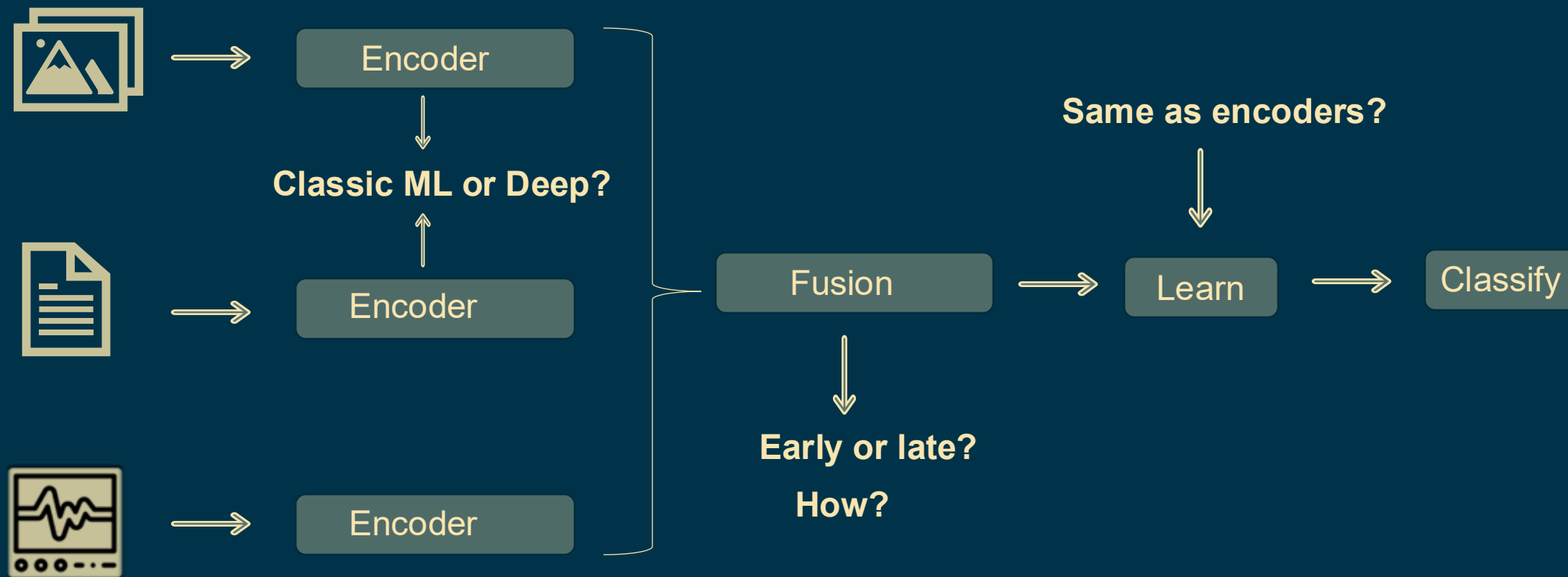
# Multimodal Chef: Recipes



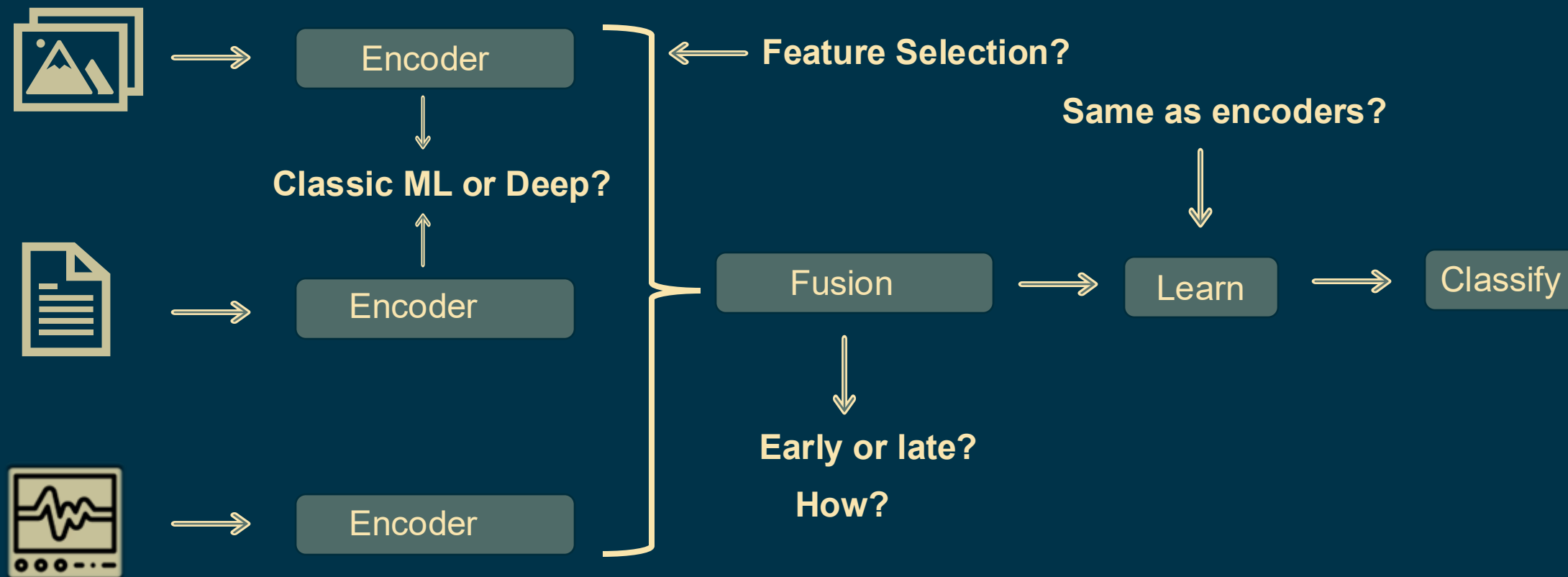
# Multimodal Chef: Recipes



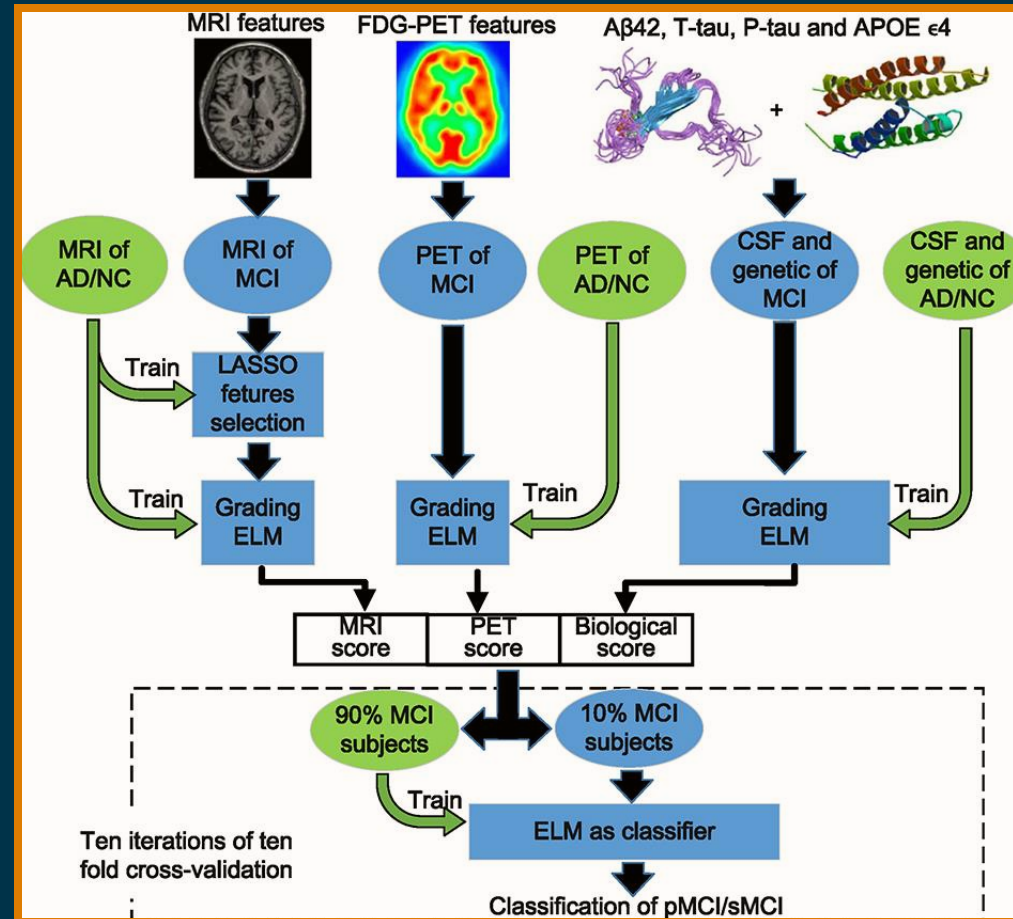
# Multimodal Chef: Recipes



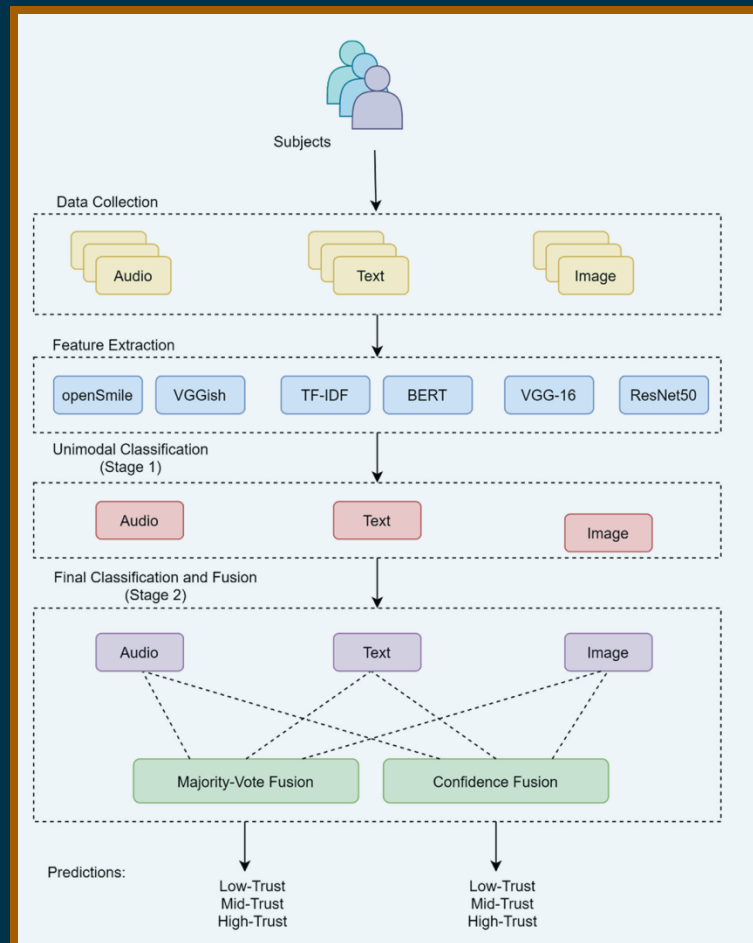
# Multimodal Chef: Recipes



# Multimodal Chef: Recipes in action

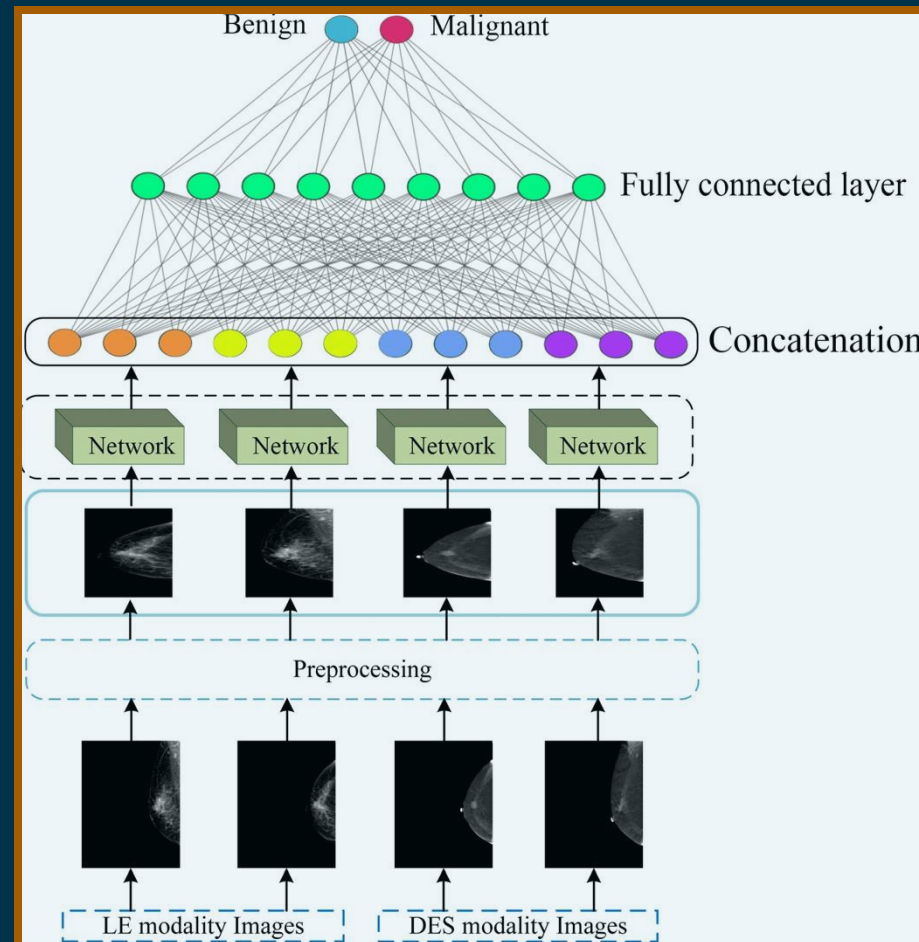


# Multimodal Chef: Recipes in action

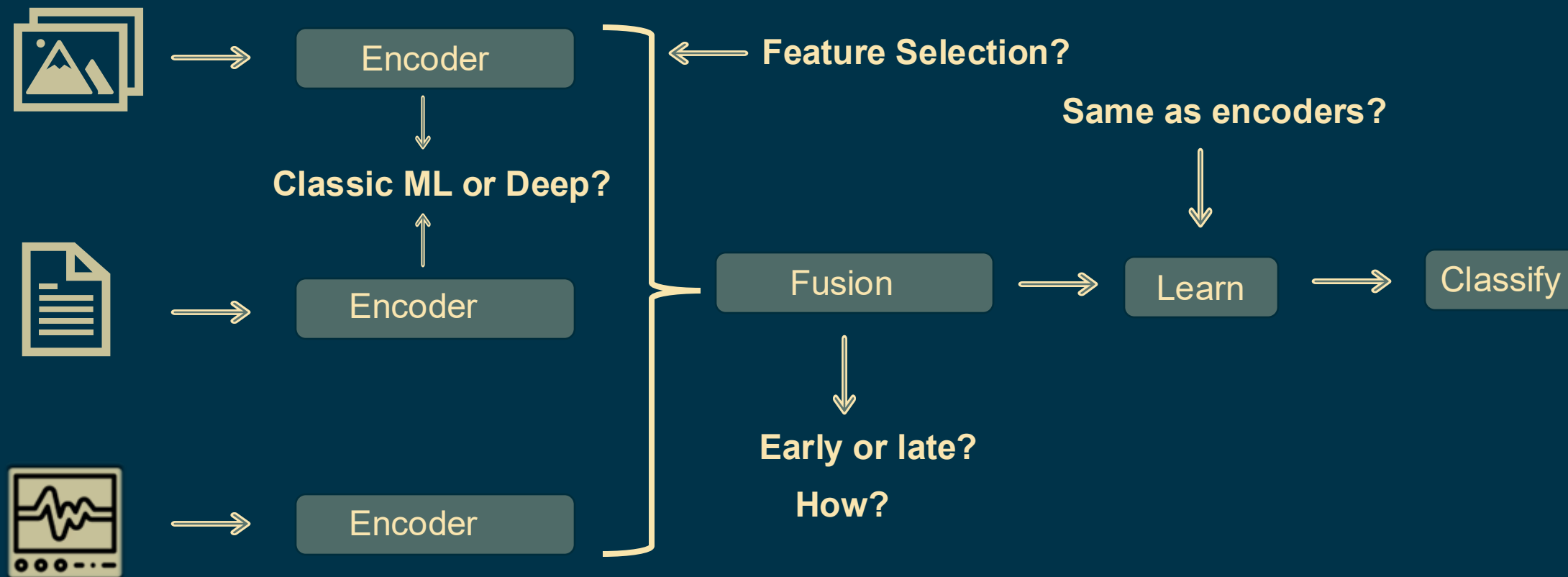




# Multimodal Chef: Recipes in action



# Multimodal Chef: Recipes



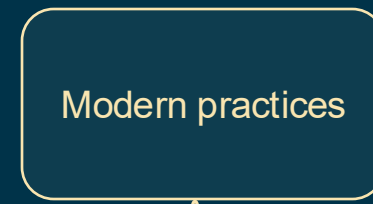
# In this talk



Motivation Challenges



Ingredients Recipes



State-of-the-art

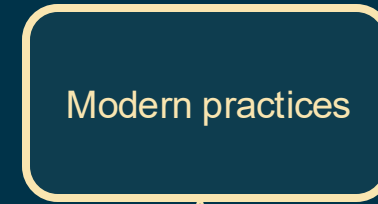
# In this talk



Motivation Challenges



Ingredients Recipes



State-of-the-art

# Recall the challenges

Representation



How do we extract relevant features?

Combination



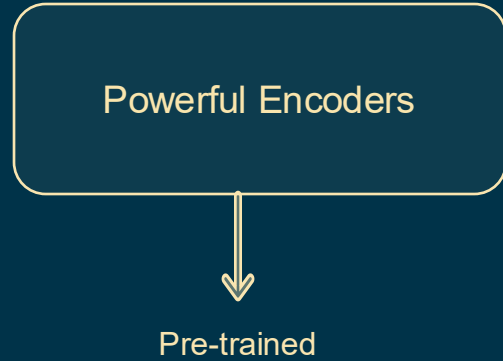
How do we fuse representations?

Labels

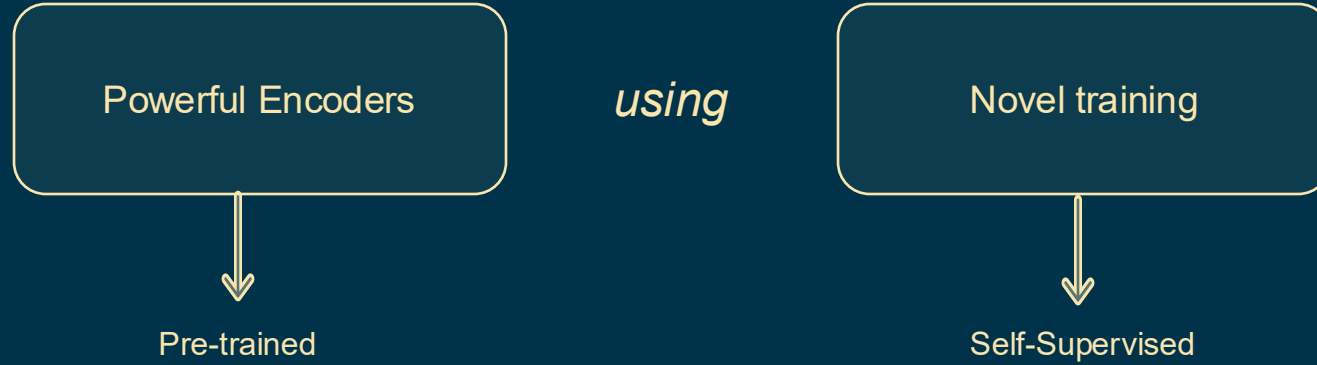


How do we learn with limited supervision?

# Modern Practice



# Modern Practice

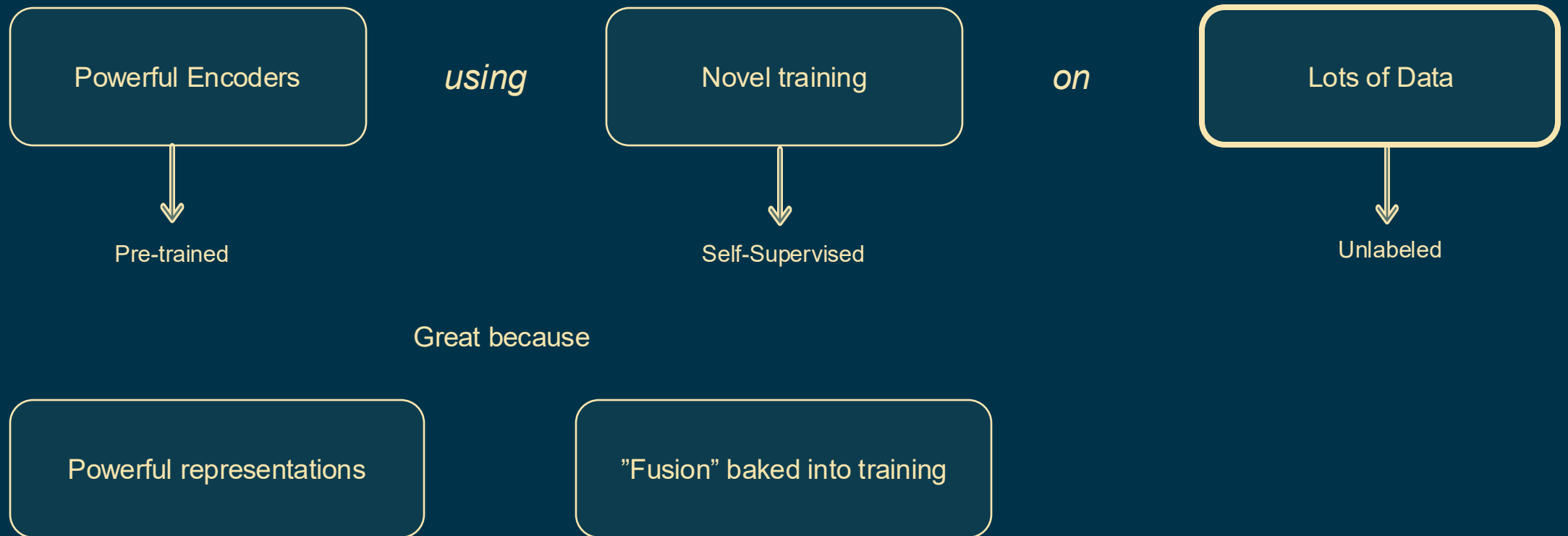


# Modern Practice

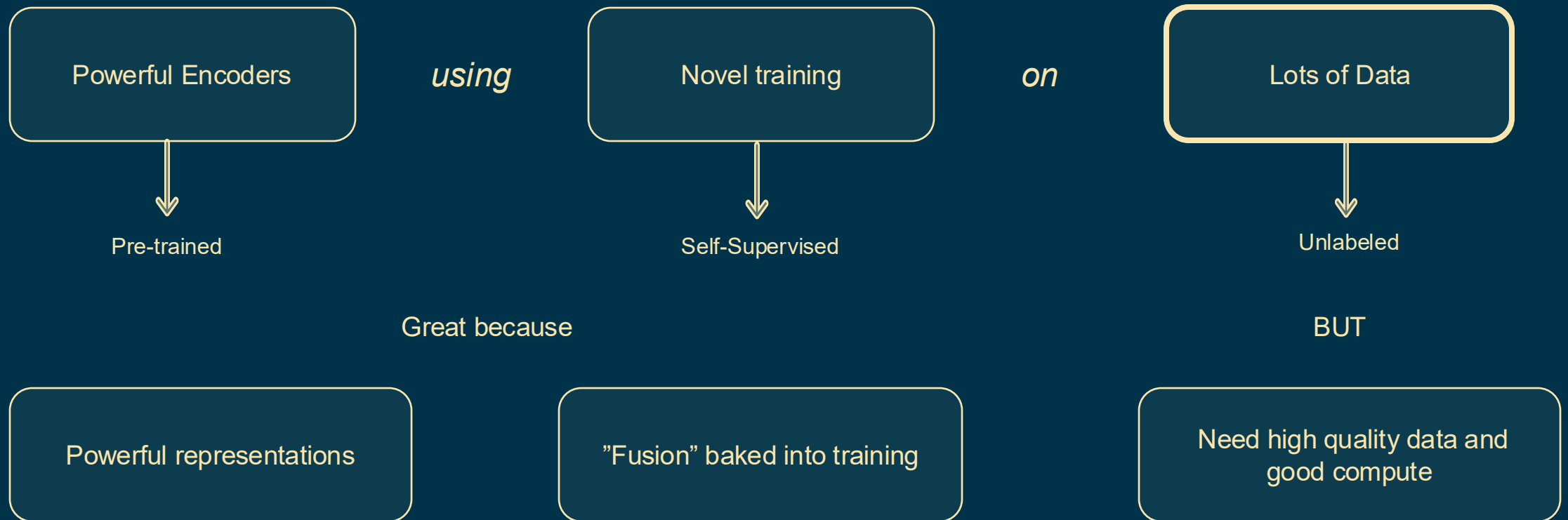




# Modern Practice

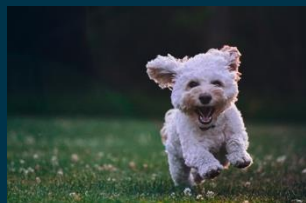


# Modern Practice



# Modern Practice: Contrastive Pre-training

Modality 1

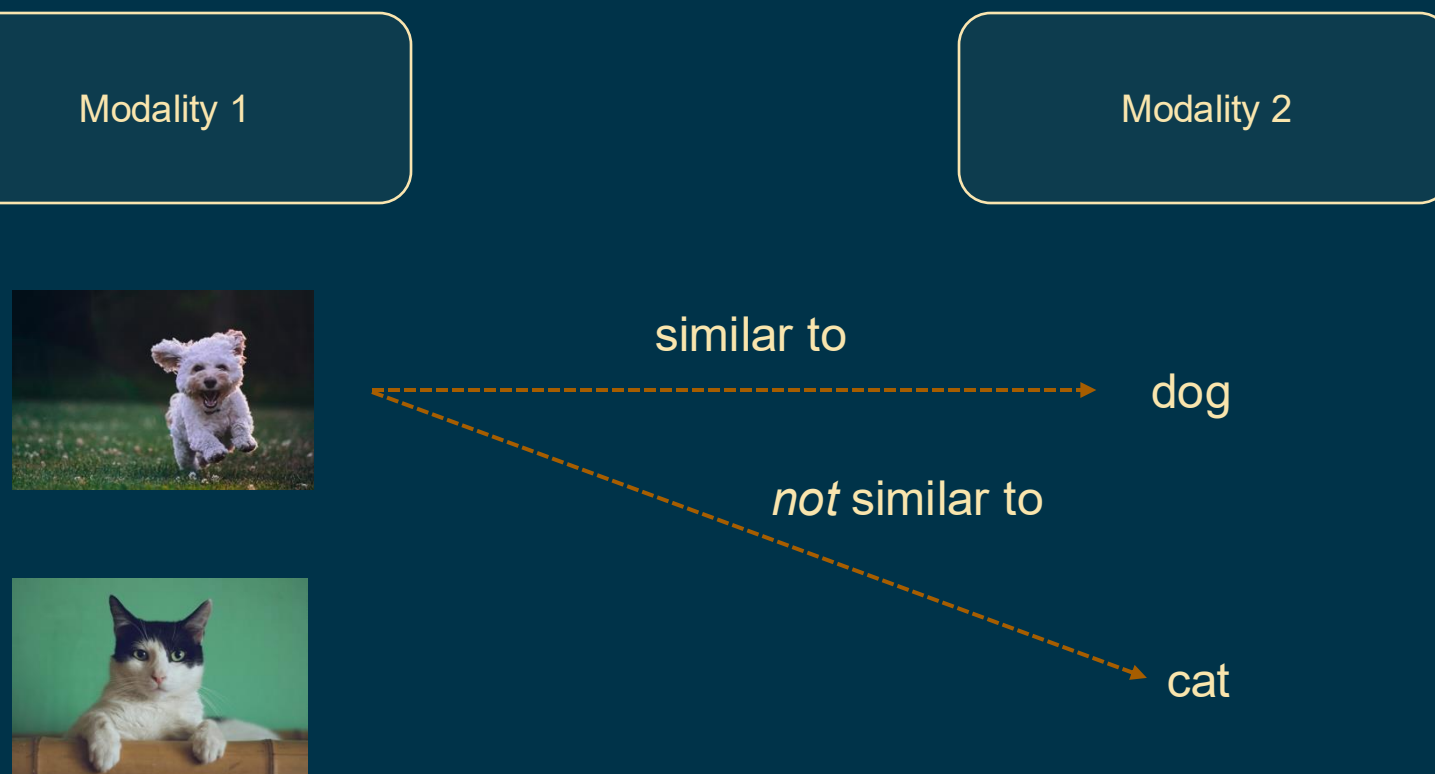


Modality 2

dog

cat

# Modern Practice: Contrastive Pre-training

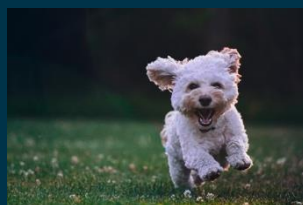


# Modern Practice: Contrastive Pre-training

You're good if you have 400 million such examples

Modality 1

Modality 2



similar to

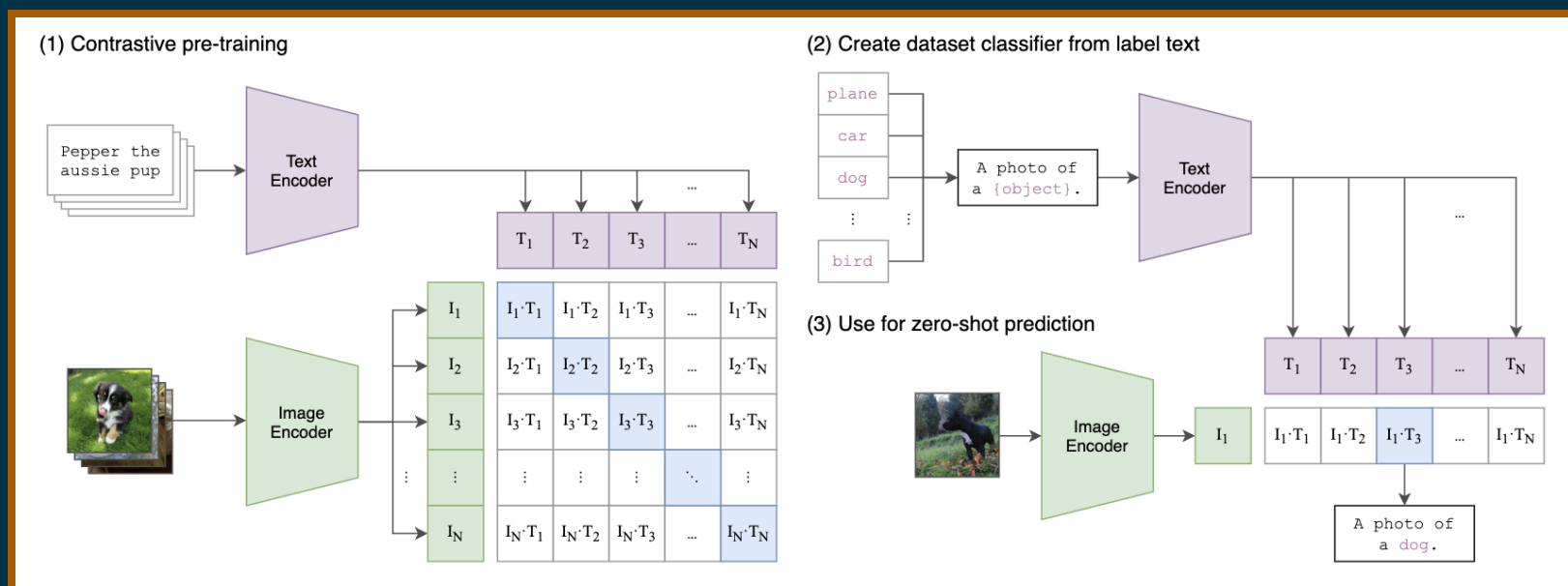
dog

*not similar to*

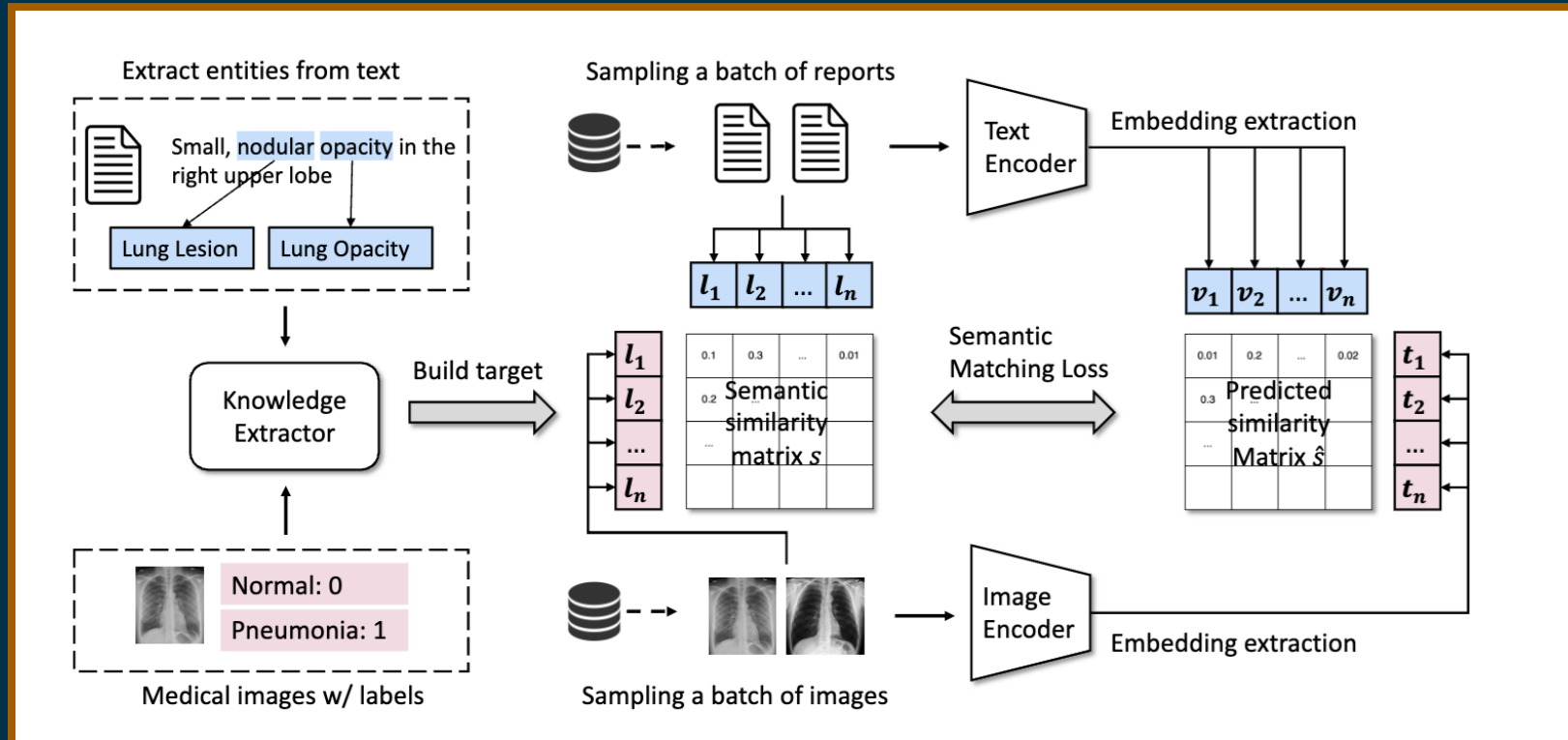
cat



# Modern Practice: Contrastive Pre-training



# Modern Practice: MedCLIP



# Modern Practice: The Contrastive Zoo!

RemoteCLIP

MedCLIP

FashionCLIP

AudioCLIP

BioCLIP

ECG-CLIP



# Modern Practice: The Contrastive Zoo!

RemoteCLIP

MedCLIP

FashionCLIP

AudioCLIP

BioCLIP

ECG-CLIP

and many more..

# Modern Practice: The Contrastive Zoo!

RemoteCLIP

MedCLIP

FashionCLIP

AudioCLIP

BioCLIP

ECG-CLIP

and many (MANY) more..

# Modern Practice: CLIP is all you need?

# Modern Practice: CLIP is all you need?

Not exactly

# Modern Practice: CLIP is all you need?

Not exactly

CLIP is simply a method of training similar representations

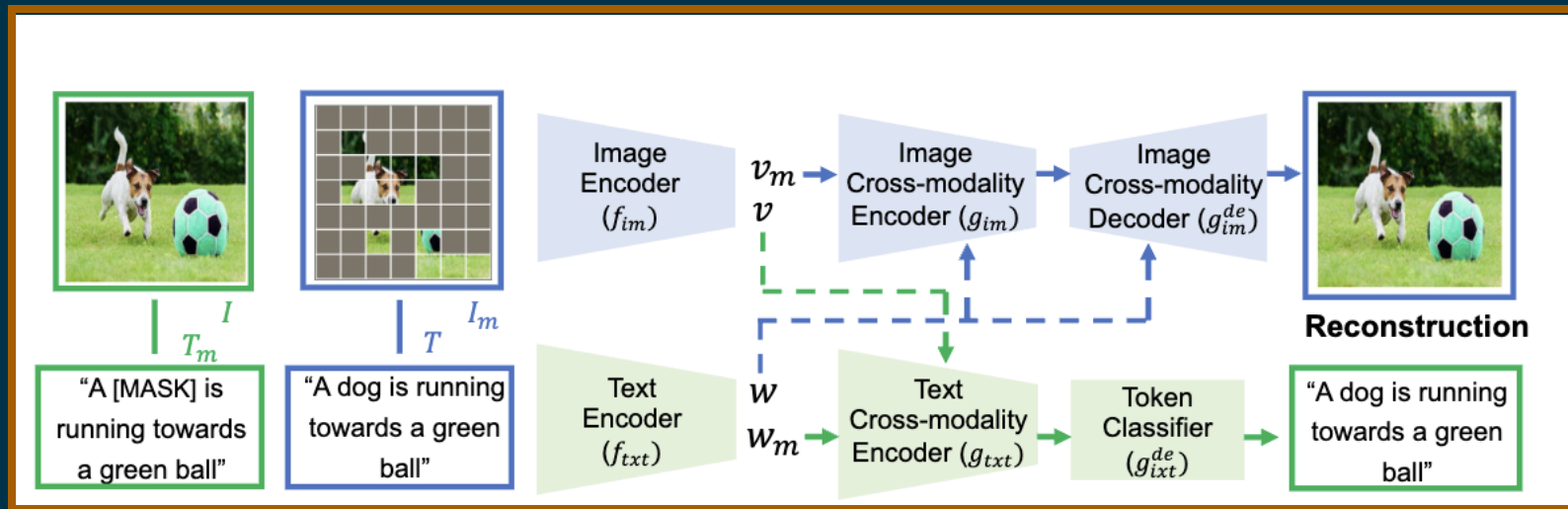
# Modern Practice: CLIP is all you need?

Not exactly

CLIP is simply a method of training similar representations

There are other quite a few popular alternatives

# Modern Practice: Masked Modelling



# Modern Practice: CLIP is all you need?

Not exactly

CLIP is simply a method of training similar representations

There are other quite a few popular alternatives



# Modern Practice: CLIP is all you need?

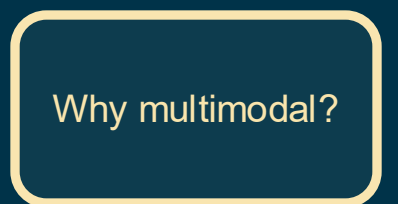
Not exactly

CLIP is simply a method of training similar representations

There are other quite a few popular alternatives

You could even pair of a powerful image-only encoder with a large language model!

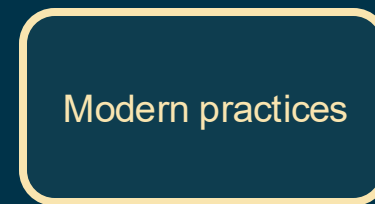
# Summary



Motivation      Challenges



Ingredients      Recipes



State-of-the-art