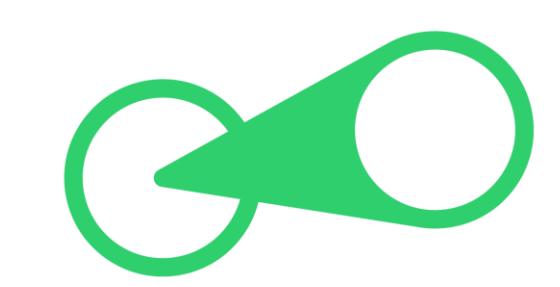


# What-Where Transformer: 物体の意味と位置情報を並行的に処理する画像バックボーンの基礎検討



Recognition, Control and Learning Algorithm Lab.



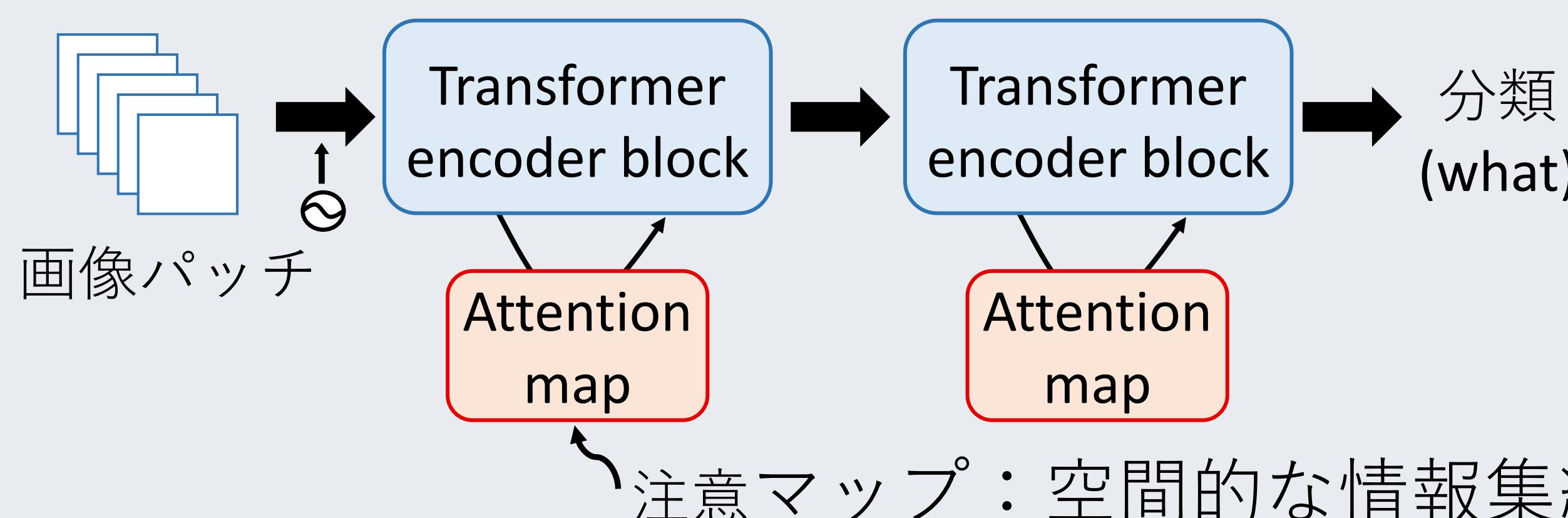
吉橋 亮太<sup>1</sup>, 加太 将弘<sup>1</sup>, 池畠 諭<sup>2,3</sup>, 川上 玲<sup>1</sup>, 佐藤 育郎<sup>1,2</sup>

<sup>1</sup>東京科学大学, <sup>2</sup>デンソーITラボラトリ, <sup>3</sup>国立情報学研究所



## ◆ 画像認識におけるTransformer

- Vision Transformer (ViT) [1]

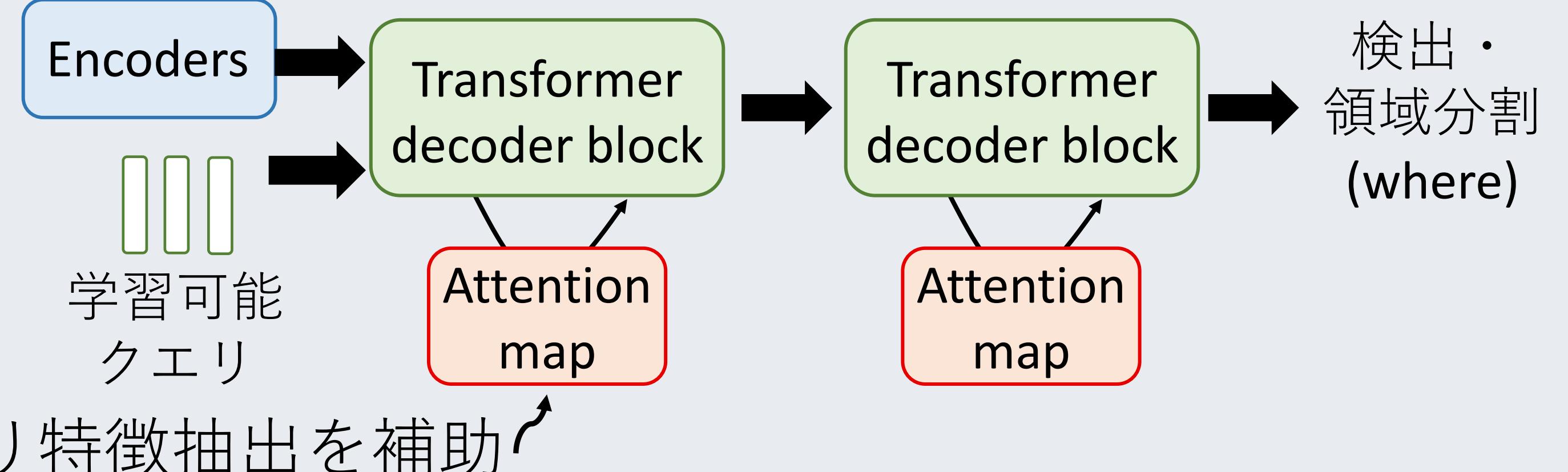


従来：分類と検出・領域分割に異なるエンコーダ・デコーダ構造が必要 = whatとwhereを別個処理

[1]: Dosovitskiy+, An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, ICLR2021  
[2]: Carion+, End-to-End Object Detection with Transformers, ECCV2020

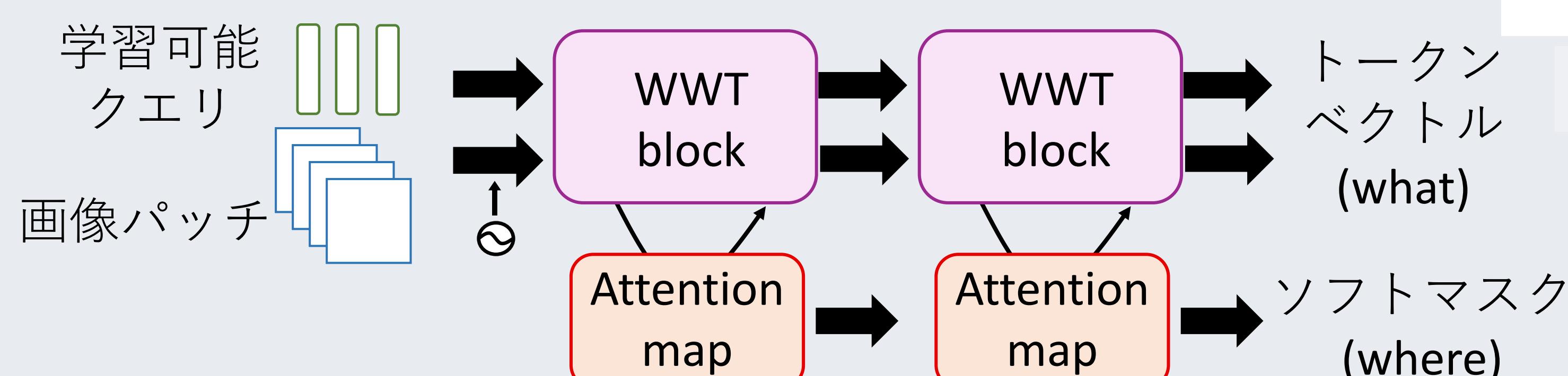
## 研究の背景と目的

- Detection Transformer (DETR) [2]



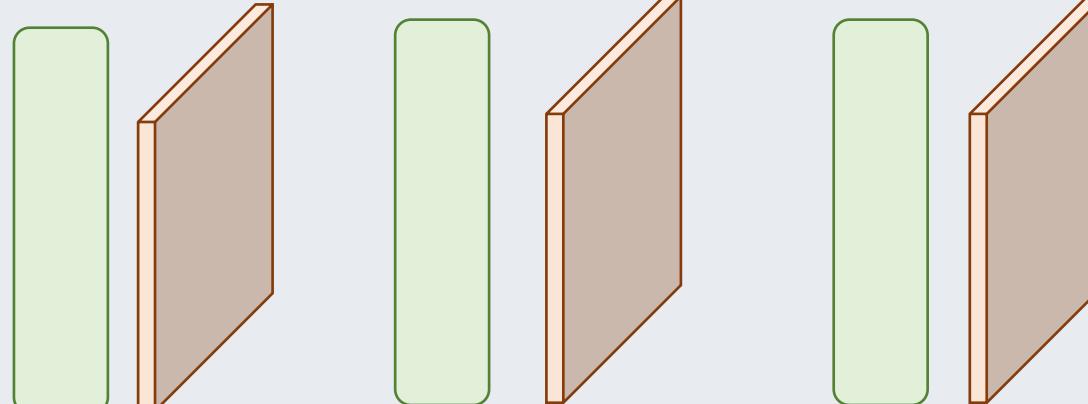
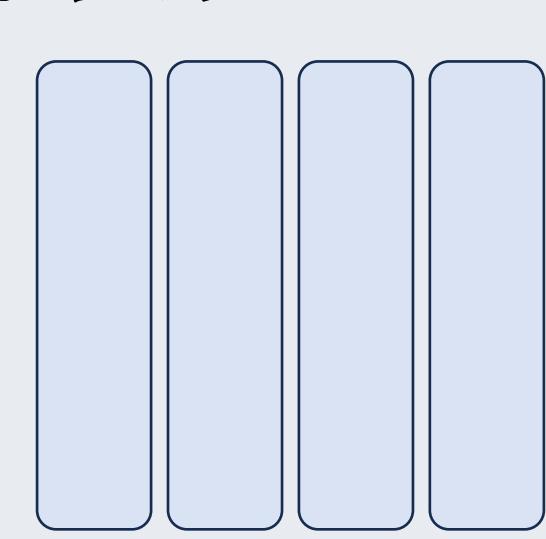
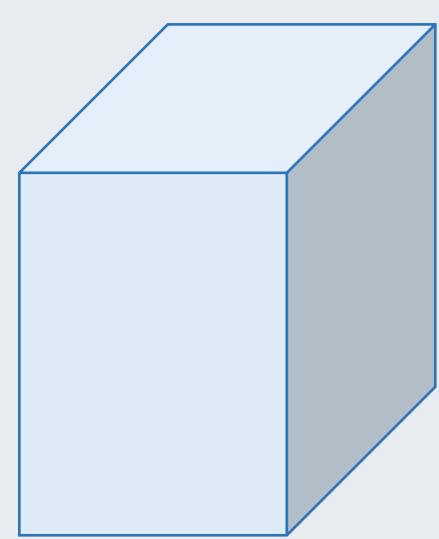
Q. 構造を統一しシンプル化・物体位置に関する情報が特徴抽出にも介在するようにできないか?  
→ ViT型バックボーン内部に物体位置情報を扱う機構を取り込む **What-Where Transformer (WWT)** を考案

## 提案手法: What-Where Transformer



- ・ トークン経路と注意マップ経路の複数ストリーム構造  
注意マップもネット内部を伝播・出力として再利用
- ・ 学習可能クエリで画像の要素（例：物体）ごとの  
トークンや注意マップを学習：この単位を”スロット”  
と呼ぶ

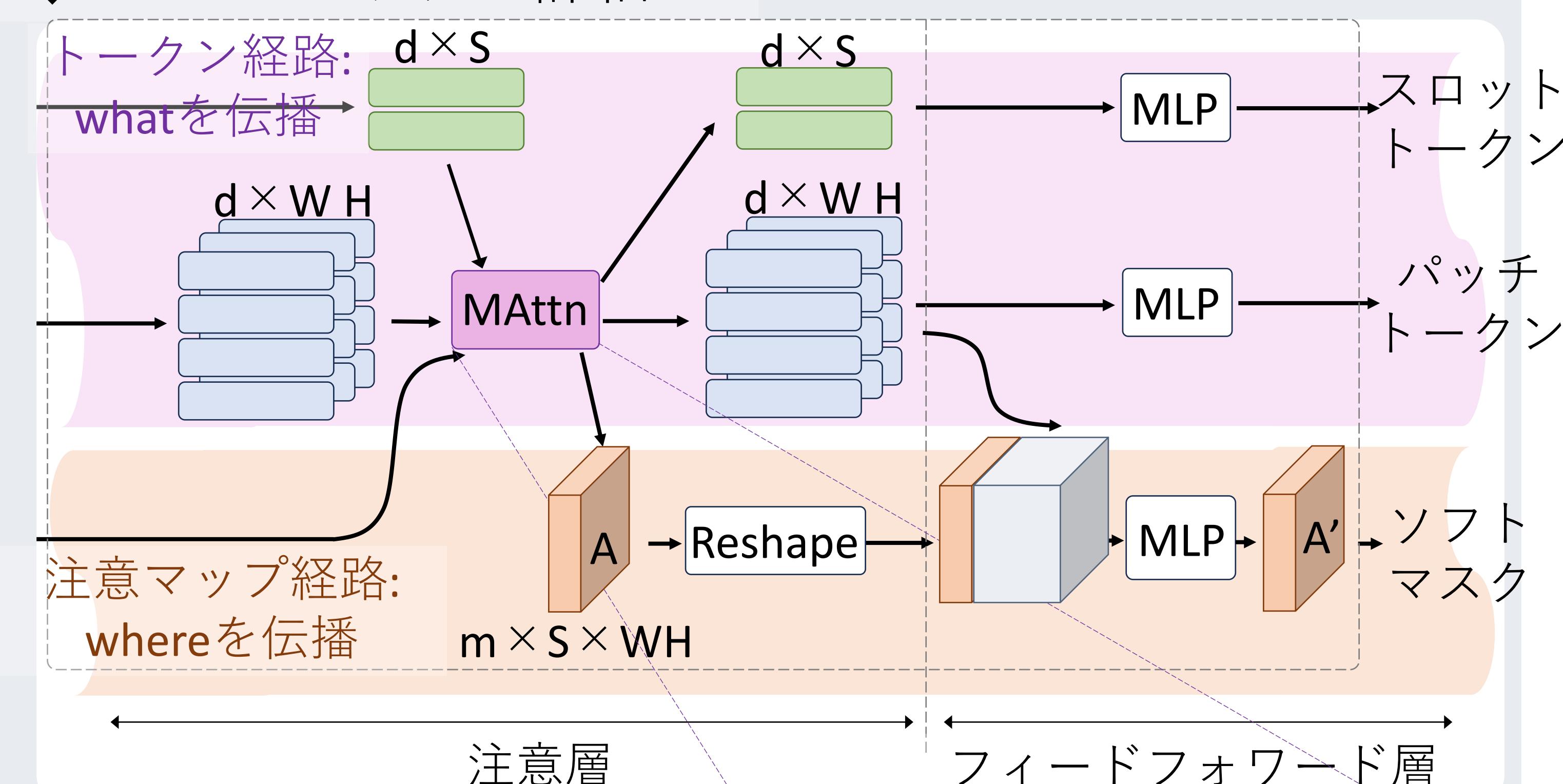
## ◆ WWTにおける画像表現



CNN: 特徴マップ  
ViT: トークン  
WWT: トークンマスク  
テンソル  
系列

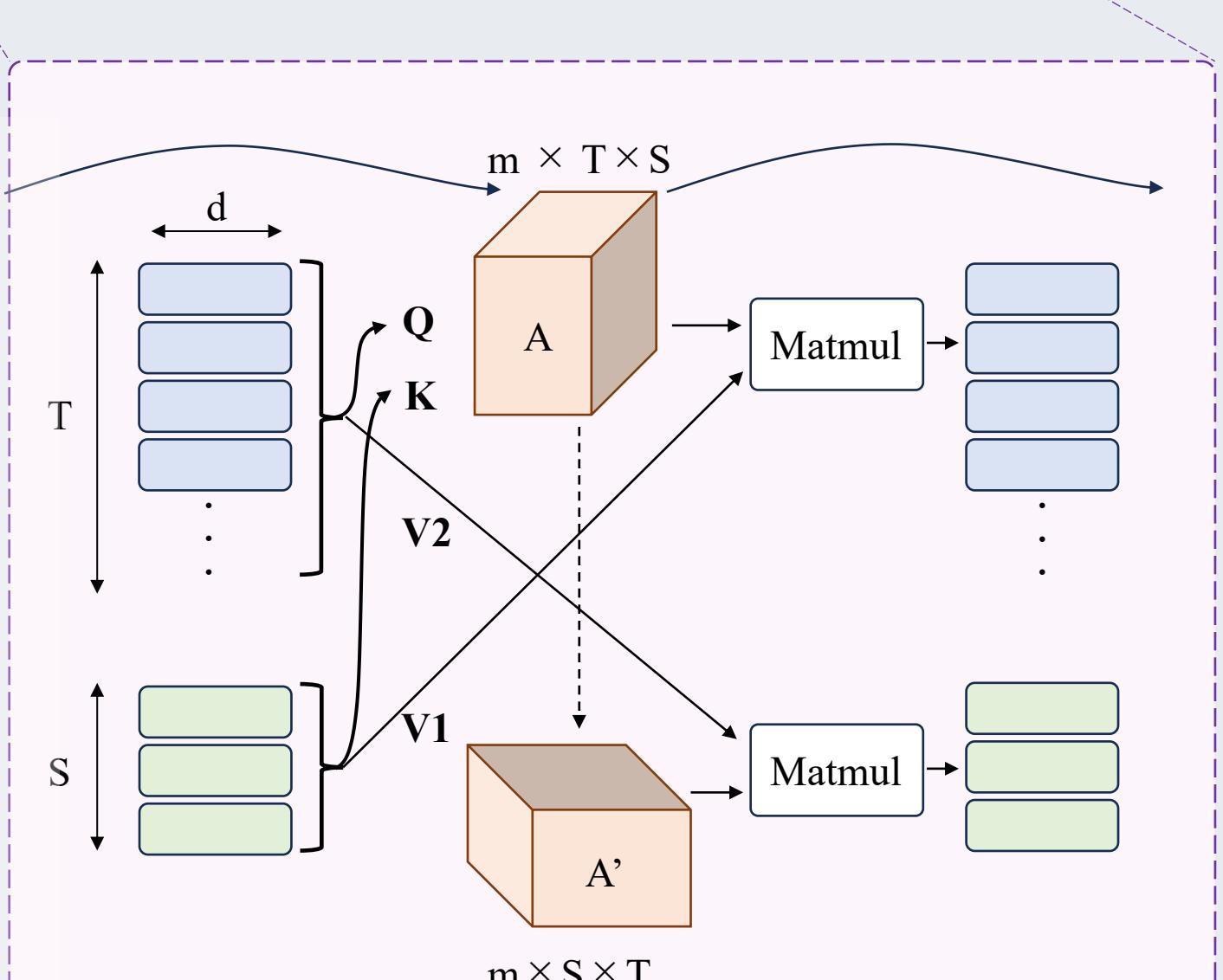
トークン-マスク  
対の集合: Whatと  
whereを分離的に  
エンコード

## ◆ WWTブロックの詳細



## ◆ 相互注意 (MAttn)

モジュール：  
パッチトークンと  
スロットトークンを  
内積型注意により  
相互作用させ  
両方を更新



## 基礎的な評価

### ◆ Multi-MNISTデータセットで

- ・ マルチラベル分類
  - ・ 入力の再構成(自己符号化)
  - ・ 教師あり領域分割
- を学習し各タスク性能と  
出力ソフトマスクを観察

### マルチラベル分類

モデル	テスト正解率 (%)	学習正解率 (%)
CNN	<b>68.94</b>	74.88
ViT	20.18	<b>99.48</b>
ViT + 相互注意	23.20	86.09
WWT (ours)	<b>32.01</b>	<b>99.61</b>

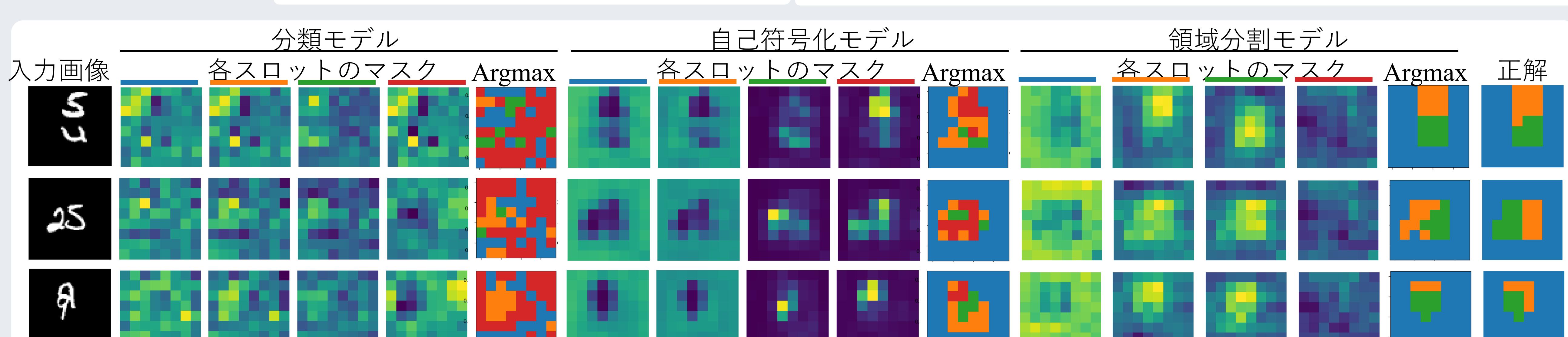
### 自己符号化

モデル	テスト mIoU	テスト再構成誤差
Slot Attention	17.20	0.0141
ViT	<b>25.56</b>	<b>0.0091</b>
WWT (ours)	<b>36.47</b>	0.0166

### 教師あり領域分割

モデル	テスト mIoU
ViT	<b>91.60</b>
WWT (ours)	88.38

◆ 出力マスクの可視化:  
自己符号化 = 教師なし学習で物体ごとのマスクをある程度獲得



◆ 今後の展望: ImageNetなどで実用的なサイズのデータセットへのスケール性を確認