

当チームの目標と研究内容

目標と背景

- 目標：橋梁などインフラ構造物の維持管理を、AI・ロボット技術で効率化・高度化
- 技術課題①：画像等を自動撮影可能なロボットの開発 ← Skydio社のドローンなど
- 技術課題②：画像等を元に**対象物の状況を理解し、説明できるAI**の開発

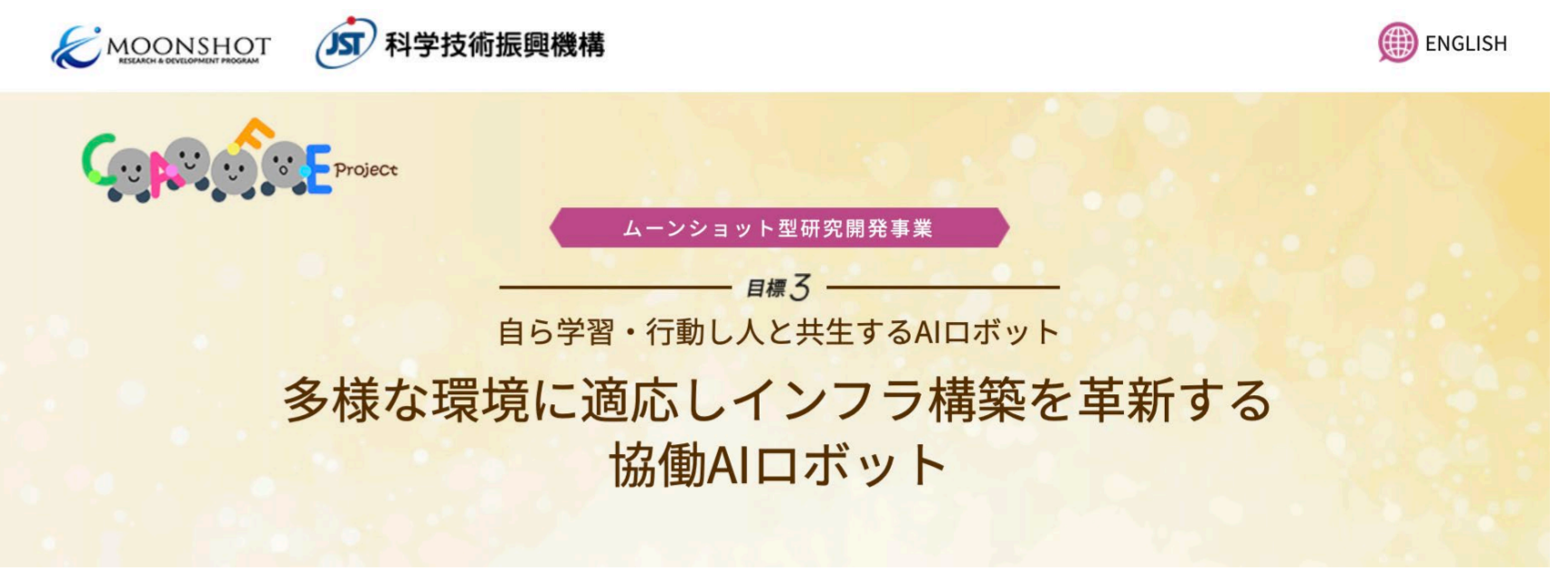


- 大規模言語モデルの発展 (ChatGPTなど) → 画像の言語化技術の重要度が増している

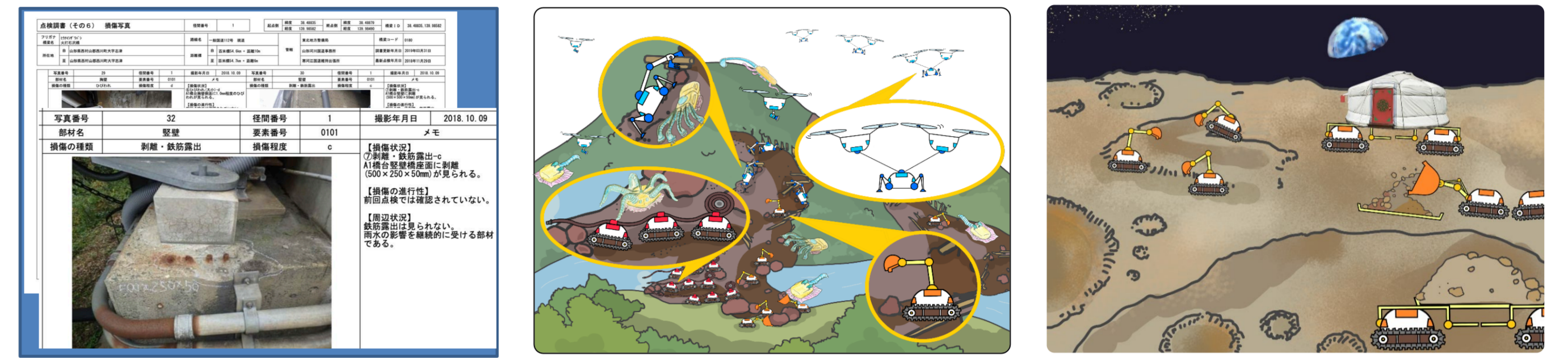
研究内容

- 画像と言語のマルチモーダルAIの研究を優先的に実施中
- 他の成果①：比喩表現を用いた画像記述 [Tran-Okatani, ACCV2022]
- 他の成果②：画像特徴の自己教師学習の研究 [Liu-Suganuma-Okatani, NeurIPS2022]

内閣府ムーンショット
CAFÉプロジェクトに
当チームが参加中



橋梁点検から災害現場・月面基地建设までをターゲットに、環境理解AIを研究中

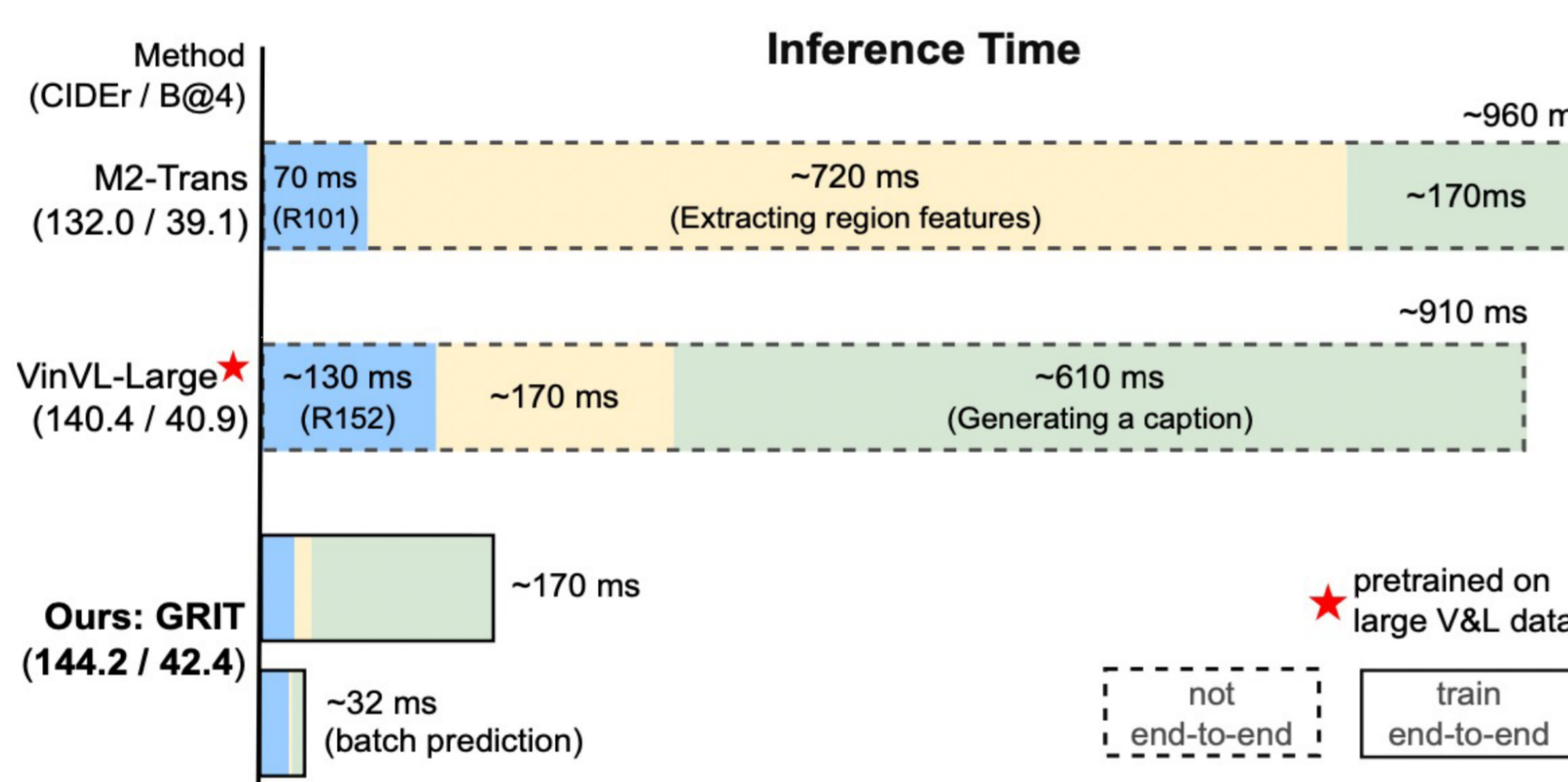


高速・高性能な画像記述AI

Nguyen, Suganuma, Okatani, GRIT: Faster and Better Image-captioning Transformer Using Dual Visual Features, ECCV2022

計算効率の高い画像記述AIの必要性

- 近年、学習が大規模化
 - モデルサイズ：順伝播計算1回~秒のオーダー
 - 訓練データサイズ：数百万~数十億の画像・テキストペアを使用
- 巨大な計算リソースがなければ、研究が不可能に
- 最適なアーキテクチャ設計で速度と性能の両方を向上



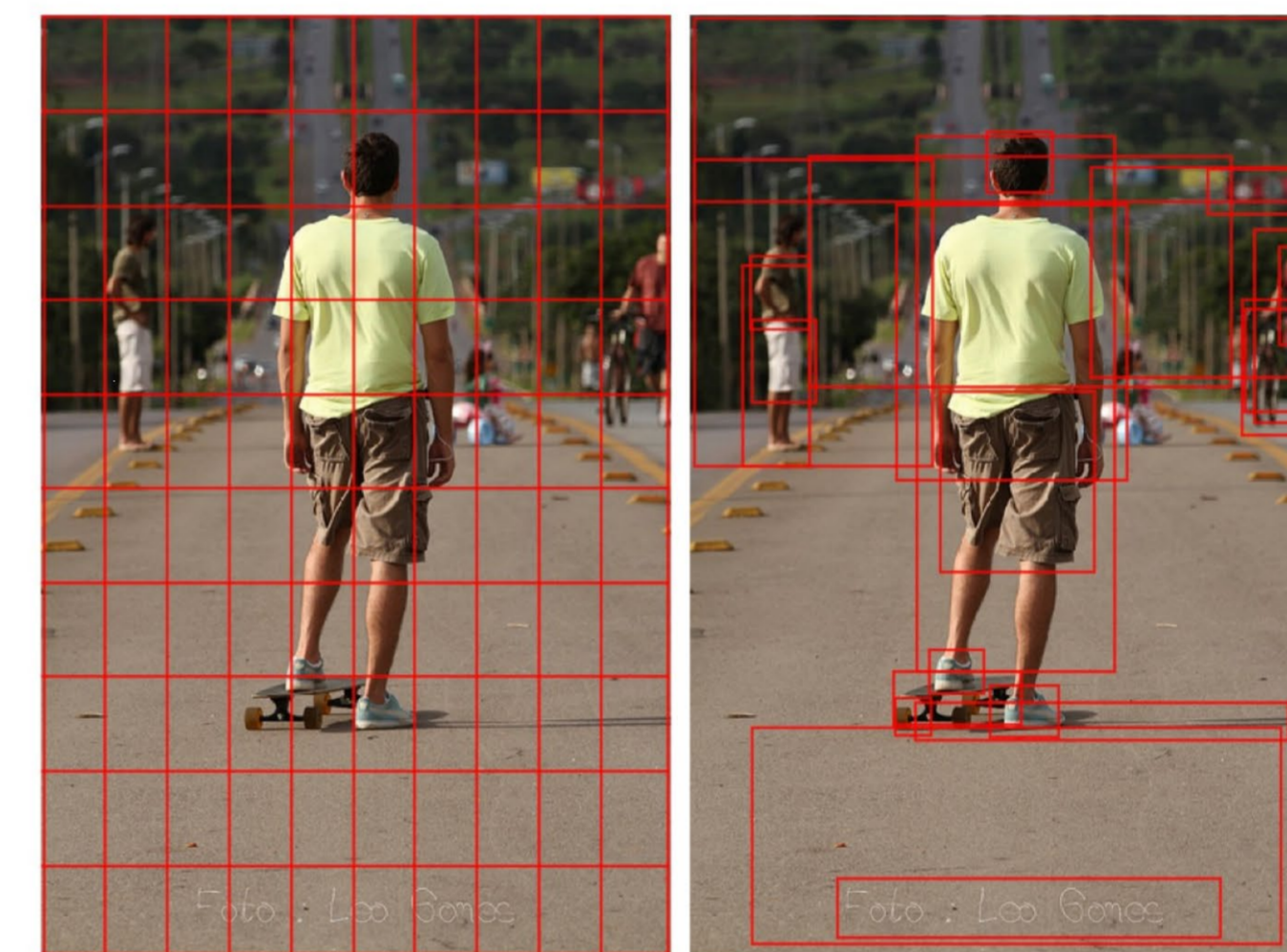
画像から取り出すべき特徴とは？

Grid-feature

- 個別物体の特徴は取れない ✗
- 物体と物体の関係性 ○
- 計算量少ない ○

Region-feature

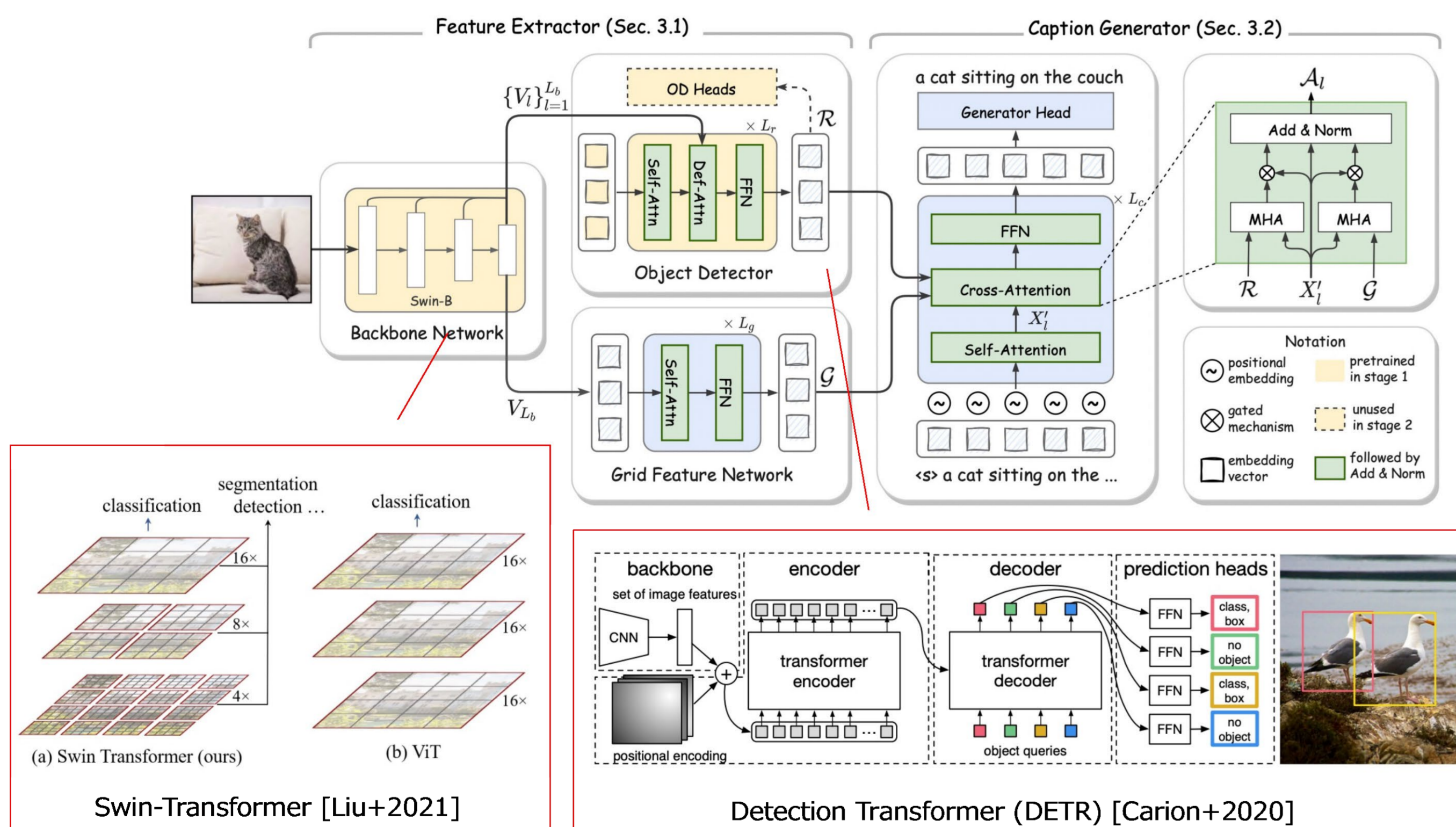
- 個別の物体の特徴を取れる ○
- 物体間関係性 ✗
- 計算量多い ✗
- 誤検出時の悪影響 ✗



[Anderson+2018]

提案した画像記述AIのアーキテクチャ (GRIT)

- Grid/Region特徴の融合利用 + "All Transformer" 構造



実験結果：記述精度

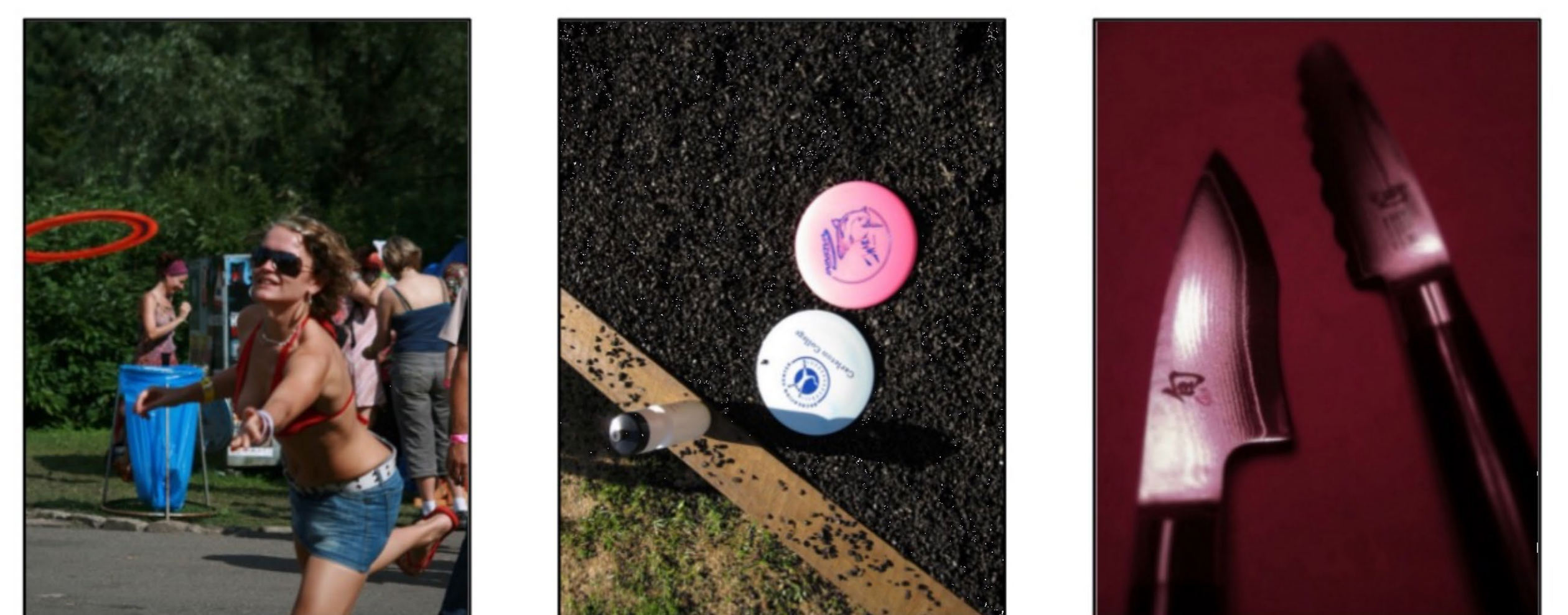
- 現時点で最良 (軽量かつ高性能)のアーキテクチャ

Method	V. E. Type	# VL Data	Performance Metrics					
			B@1	B@4	M	R	C	S
w/ VL pretraining								
UVLP [55]	R	3.0M	-	39.5	29.3	-	129.3	23.2
Oscar _{base} [25]	R	6.5M	-	40.5	29.7	-	137.6	22.8
VinVL [53]	R	8.9M	-	41.0	31.1	-	140.9	25.2
SimVLM _{large} [45]	G	1.8B	-	40.6	33.7	-	143.3	25.4
w/o VL pretraining								
SAT [43]	G	-	-	31.9	25.5	54.3	106.3	-
SCST [38]	G	-	-	34.2	26.7	55.7	114.0	-
LSTM-A [52]	G	-	-	78.6	35.5	27.3	56.8	118.3
RSTNet [54]	G	-	-	81.8	40.1	29.8	59.5	135.6
Up-Down [4]	R	-	-	79.8	36.3	27.7	56.9	120.1
RFNet [20]	R	-	-	79.1	36.5	27.7	57.3	121.9
GCN-LSTM [51]	R	-	-	80.5	38.2	28.5	58.3	127.6
LBPF [35]	R	-	-	80.5	38.3	28.5	58.4	127.6
SGAE [49]	R	-	-	80.8	38.4	28.4	58.6	127.8
AoA [15]	R	-	-	80.2	38.9	29.2	58.8	129.8
GET [16]	R	-	-	81.5	39.5	29.3	58.9	131.6
ORT [13]	R	-	-	80.5	38.6	28.7	58.4	128.3
ETA [24]	R	-	-	81.5	39.3	28.8	58.9	126.6
M ² -Transformer [7]	R	-	-	80.8	39.1	29.2	58.6	131.2
X-LAN [33]	R	-	-	80.8	39.5	29.5	59.2	132.0
TCIC [10]	R	-	-	81.8	40.8	29.5	59.2	135.4
Dual Global [46]	R+G	-	-	81.3	40.3	29.2	59.4	132.4
DICT [31]	R+G	-	-	81.4	39.8	29.5	59.1	133.8
GRIT	R+G	-	-	83.5	41.9	30.5	60.5	142.2
GRIT [†]	R+G	-	-	84.2	42.4	30.6	60.7	144.2

画像記述の例



- | | | | |
|-------------------------|--|--|---|
| 正解 (人間)
既存手法
提案手法 | GT-1: An elderly man looks at a cell phone.
GT-2: An old man holding up a cell phone to his face. | GT-1: a table top with some plates of food on it
GT-2: Two plates of breakfast foods on a restaurant table. | GT-1: there are many people in the beach playing volleyball
GT-2: some males on some sand are playing volleyball |
| | M ² : a man is taking a picture of himself on a motorcycle | M ² : a plate of food with eggs and meat on a table | M ² : a group of people playing soccer on the beach |
| | GRIT: a man sitting in a chair holding a cell phone | GRIT: two plates of food on a table with a fork | GRIT: a group of men playing volleyball on the beach |



- | | | | |
|--|---|---|--|
| 正解 (人間)
既存手法 (M2-transformer)
提案手法 | GT-1: A woman throwing a frisbee outside at a park
GT-2: a woman is throwing a disk outside in the sun | GT-1: Two frisbees laying on the ground next to a sports water bottle.
GT-2: Two flying disks on the ground next to a water bottle | GT-1: Two knives are lying on a dark red surface.
GT-2: Two knives placed on a dining table |
| | M ² : a woman holding a blue umbrella in the street
GRIT: a woman is throwing a frisbee in the street | M ² : a knife and a knife on a table with a | M ² : a close up of a red tie with a |
| | GRIT: two frisbees laying on the ground next to a bottle | GRIT: two knives are on a red table with | |