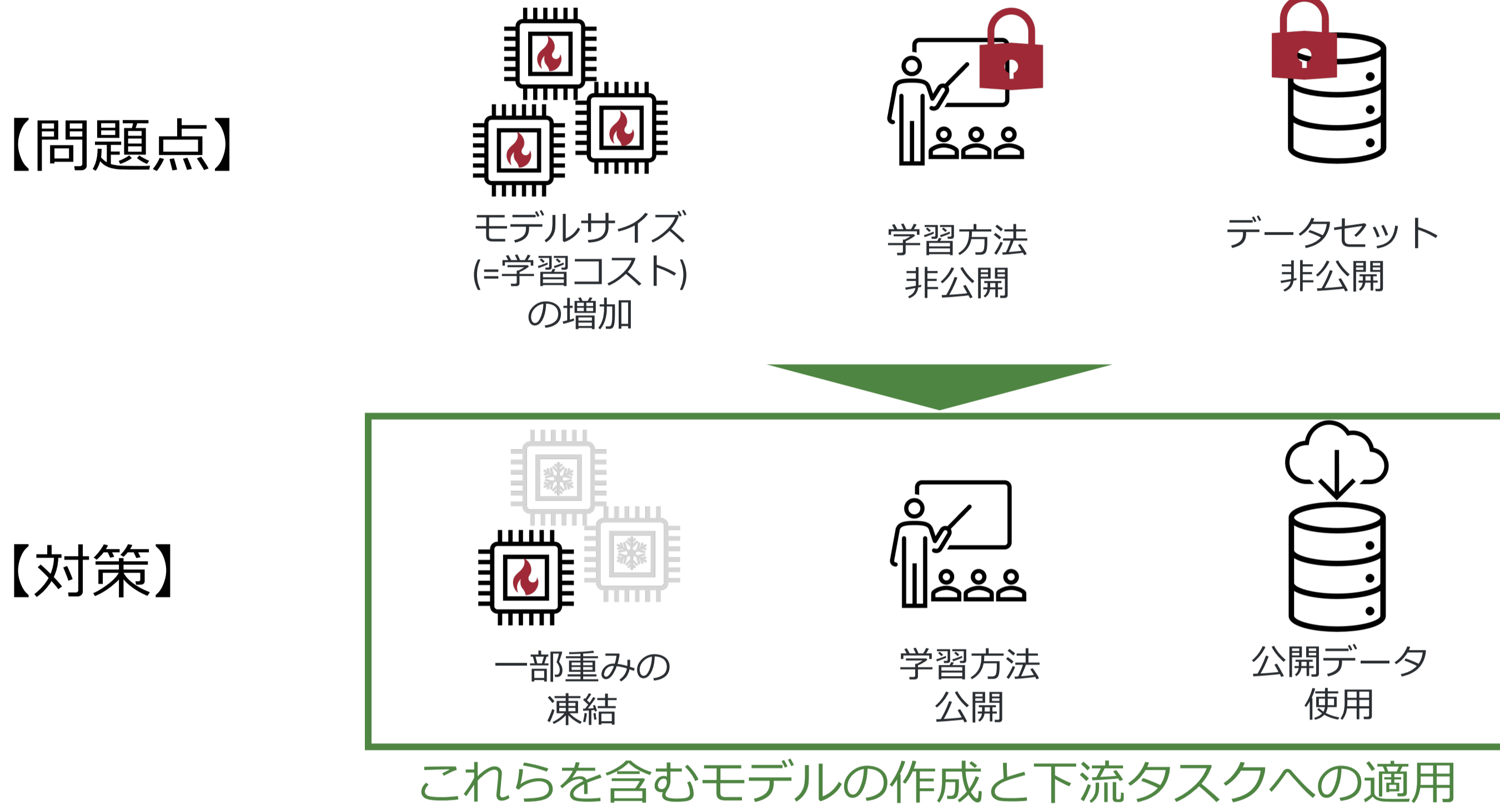


# 大規模Vision&Languageマルチタスクモデルの学習効率化と Human-Object Interactionへの適用

軸屋敬介<sup>+</sup>, 梁瀬和哉<sup>+</sup>, 表英輝<sup>+</sup>, 土田裕登<sup>+</sup>, 加藤邦人<sup>+</sup>  
<sup>+</sup>: 岐阜大学工学部, <sup>+</sup>: 日本車輛製造株式会社

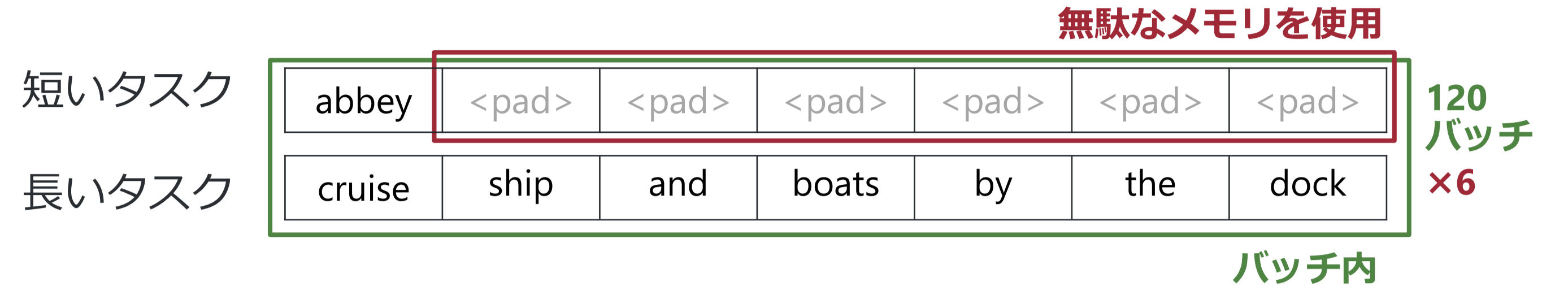
## 研究背景

- ◆ 大規模Vision&Languageモデル  
 画像と言語の統合 → 高い性能
- ◆ 大規模マルチタスクモデル  
 知識の組み合わせ → 未学習タスクを解く (推論能力)

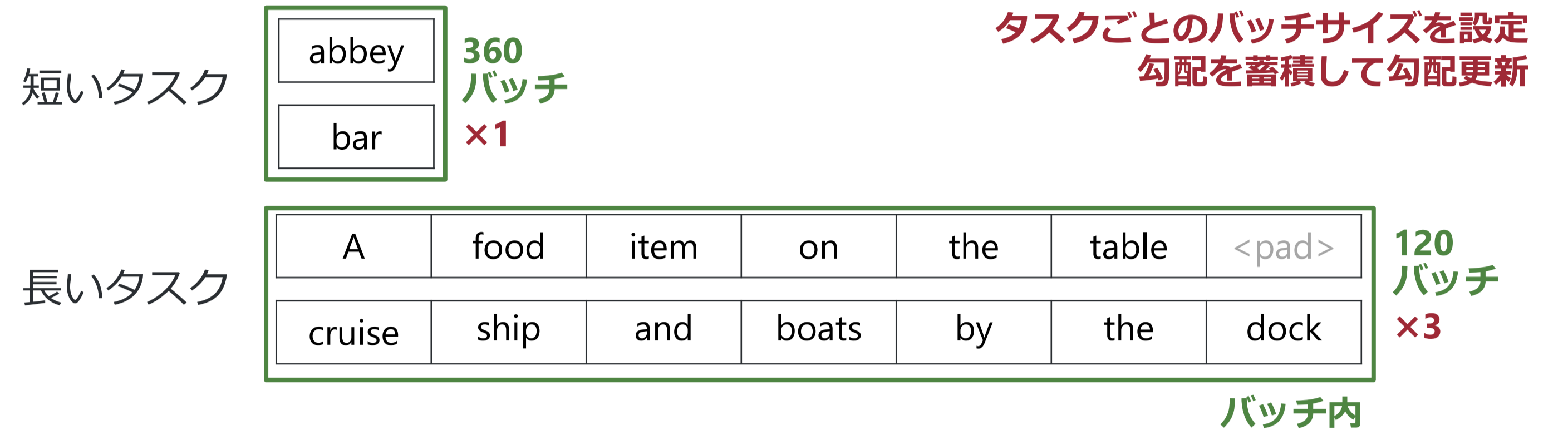


## 【効率的な学習】

◆ 720データ (短い:長い=360:360) での計算  
**通常マルチタスク学習**



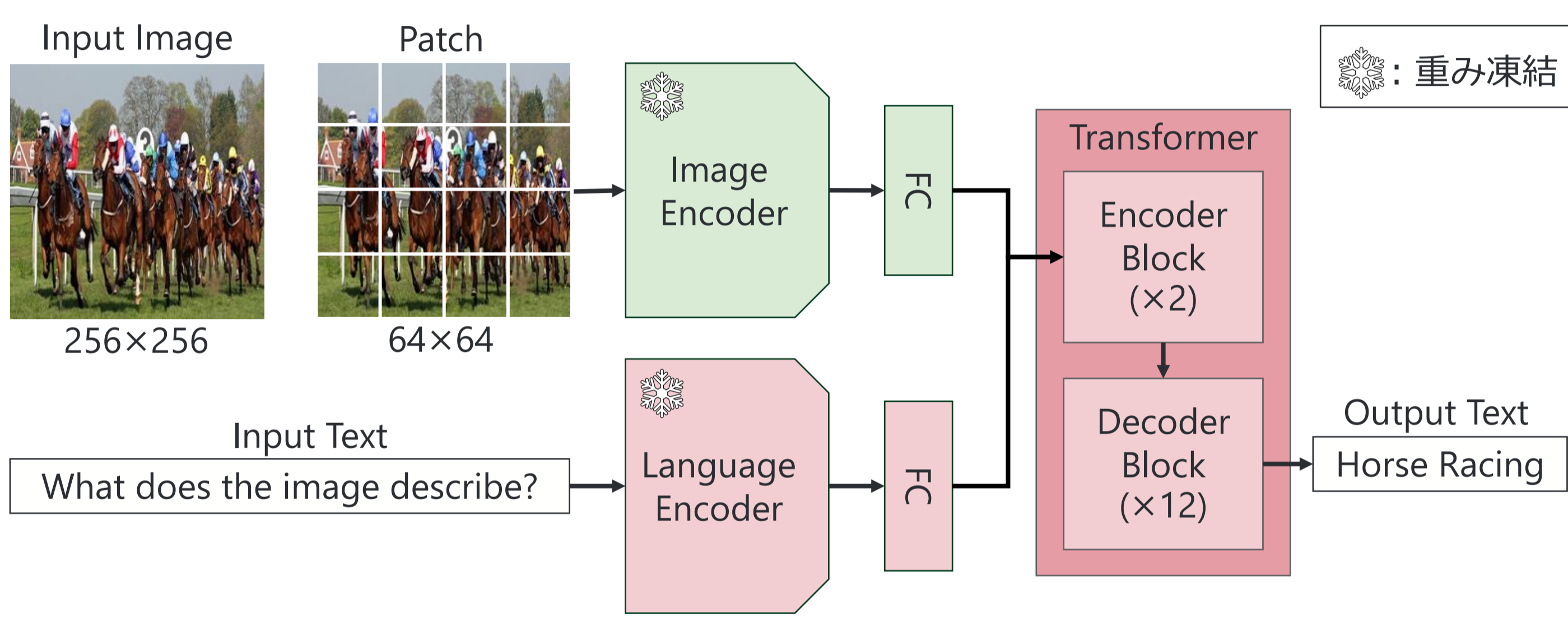
## 提案手法



→ 6 から 4 (=1+3) 回へ **計算回数の削減**

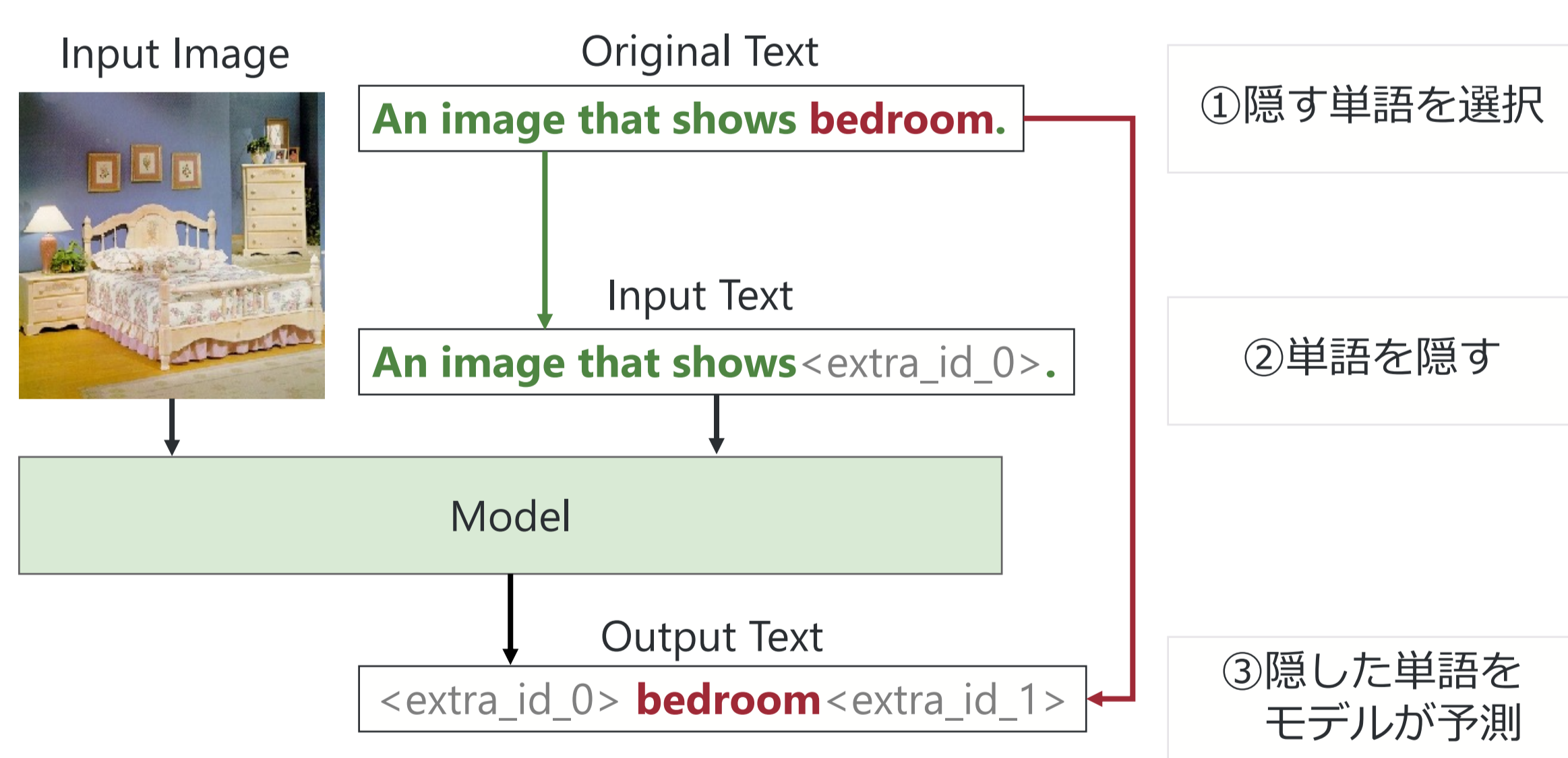
## 提案手法

### 【モデル図】



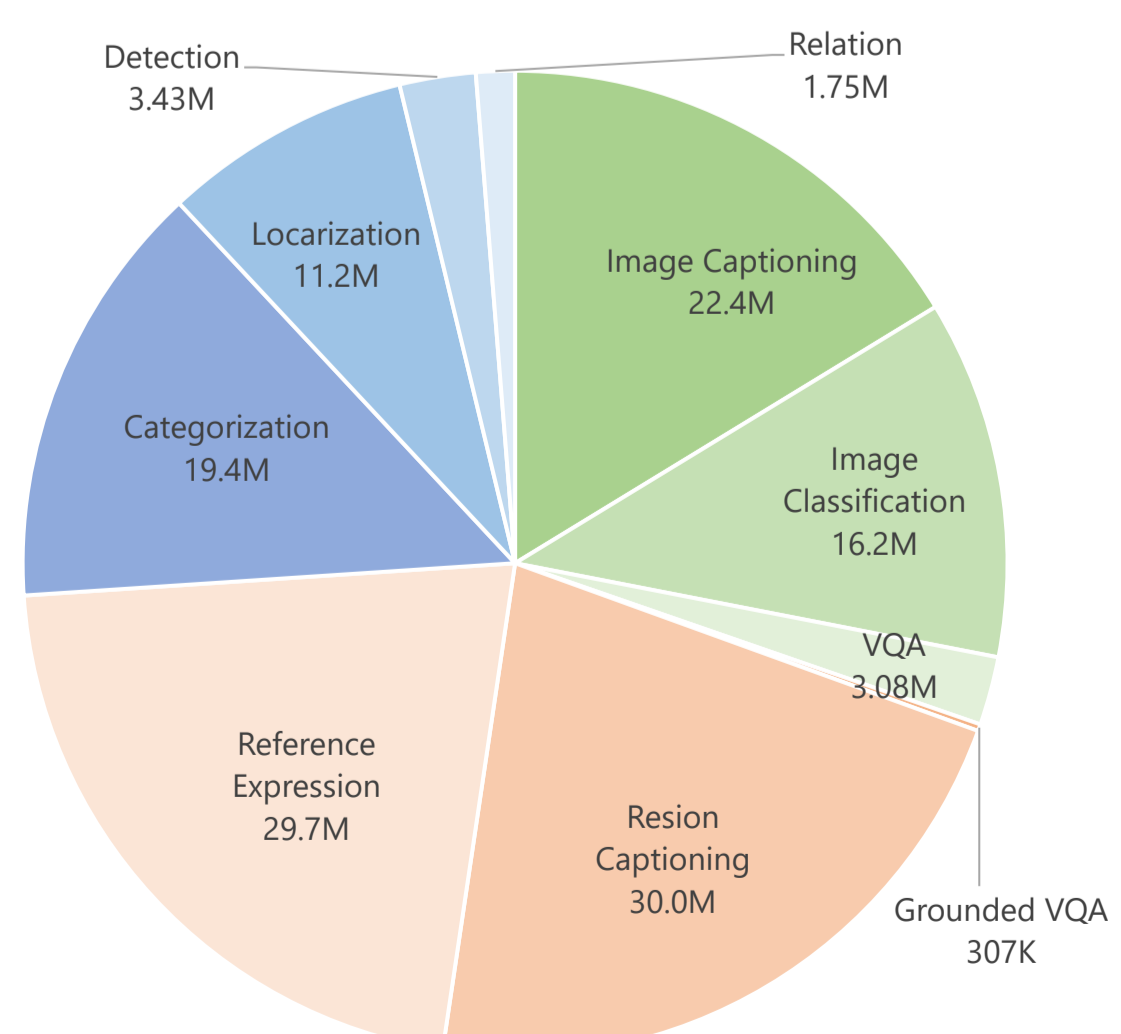
### 【事前学習】

文章の一部を隠して学習



### 【タスク学習】

様々なタスクで汎用的な知識を得る



### 【ファインチューニング】

特定のタスクに適用

**Human-Object Interaction**  
 人と物体の関係を予測するタスク



人がボールを持っている

タスクの例

タスク	入力	出力
Region Captioning	What does the region <loc_81><loc_1192> describe?	heavily embellished wedding sandals
Categorization	What is the category of the region <loc_825><loc_1349>?	Human arm

## 実験結果

### データセット

V-COCO

データセットの詳細		
学習	テスト	関係
8,543	7,811	29

### 入力と出力

入力: What is the interaction between person<loc\_40><loc\_1558> and donut<loc\_376><loc\_1024>?  
 出力: hold, eat

### 評価指標

Marco F1

$$\text{Macro F1} = \frac{1}{m} \sum_{i=1}^m (F1)_i$$

### 比較対象

OFA, Kosmos-2

### 結果

vs OFA<sub>Large</sub>  
 11倍 Δ 13.3% ↓

vs OFA<sub>Huge</sub>  
 24倍 Δ 11.5% ↓

vs Kosmos-2  
 120倍 Δ

結果の例

画像と領域	正解	OFA <sub>Large</sub>	OFA <sub>Huge</sub>	Ours	Kosmos-2
	look, kick	look	look, kick	look, kick	look, kick
	lay	sit	sit	lay	lay

パラメータと学習時間およびMacro F1

モデル	OFA <sub>Large</sub>	OFA <sub>Huge</sub>	Ours	Kosmos-2
総パラメータ	473M	929M	699M	1798M
学習パラメータ	412M	853M	162M	1798M
学習時間 [1エポック] (6000Ada×4)	11 min	24 min	1 min	120 min
RTX4090でのバッチサイズ	16	不可	160	不可
Macro F1 (Oursとの比較)	68.0% (-13.3%)	69.8% (-11.5%)	81.3% (±0%)	84.3% (+3.0%)

## 今後の課題

パラメータ数が最多のKosmos-2が最も良い精度

- モデルやデータの改良で軽量を維持しつつKosmos-2を超える性能を実現したい