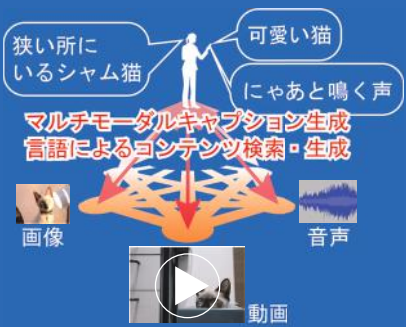


# あなたの見たものを言葉にする技術

## 未来ビジョン

- 言語をクエリとしたメディアの検索・生成
- 視聴に時間のかかるメディアをキャプション生成して表示

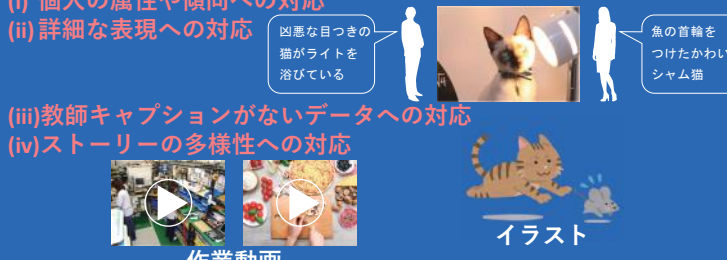


多様なデータへのキャプション生成



## 未来ビジョン達成のために必須+未達の4要求機能

- (i) 個人の属性や傾向への対応
- (ii) 詳細な表現への対応
- (iii) 教師キャプションがないデータへの対応
- (iv) ストーリーの多様性への対応



## 個人の属性や傾向への対応

### 個人へのキャプション生成転移に向けた現実的なドメイン適応

ドメイン適応：教師付きデータ（ソースドメイン）から教師が無い/欠けたデータ（ターゲットドメイン）への転移

ただし通常のドメイン適応はラベルノイズが無く、ソースとターゲットでクラスが一致している → ラベルノイズがあり、ソースとターゲットで潜在的にクラスが一致していないドメイン適応

### 2つの識別器による Noisy Universal Domain Adaptation

#### A) ノイズ除去ステップ:

- ソースデータでの学習  
識別器間の意見の食い違いと推定の曖昧さ+損失が大きい=ノイズ  
→ 上記のようなサンプルを除外した上でモデルの初期値を学習
- ターゲットデータでの学習  
2つの識別器の食い違いまたは識別器での曖昧さが一定以下  
→ ソースにも存在するクラス  
一定以上 → ターゲットにしか存在しないクラス

#### B) クロスエントロピー最大化ステップ:

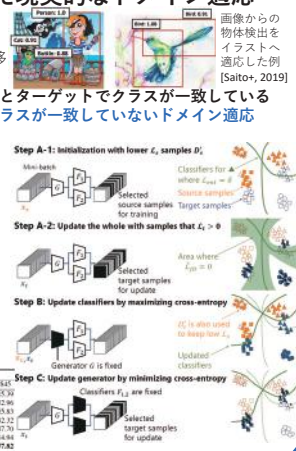
2つの識別器の食い違いを最大化するように識別器を更新

#### C) クロスエントロピー最小化ステップ:

2つの識別器の食い違いを最小化するように分布（特徴量）を更新

上記Step BとStep Cをミニバッチごとに繰り返す

実験結果  
提案手法が最もノイズ(S/P<sub>no</sub>)に強いことを複数データで確認



CVPR 2021 投稿中

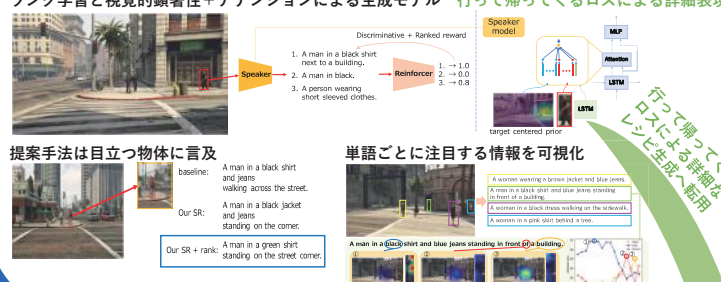
## 詳細な表現への対応

### 人がすばやく見分けられる差分に注目した記述生成

#### CGデータセットの収集+公開

CGデータセットの収集+公開  
ランク学習と視覚的顕著性+アテンションによる生成モデル

ランク学習と視覚的顕著性+アテンションによる生成モデル  
行って帰ってくるロスによる詳細表現



ICCV 2019 採録!

## 教師キャプションがないデータへの対応

### 疑似キャプションを用いた教師なしキャプション生成

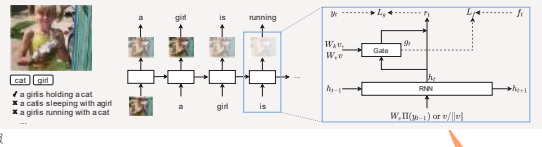
画像と文がバラバラのデータ → 含まれる物体が同じ画像と文を疑似的にペア化・学習

#### ゲートの活用

- ✓ 画像を用いた単語推定
- ✓ 生成中の文を用いた単語推定

#### ゲートの疑似教師

- ✓ ゲートの教師情報は存在しない
- ✓ だが疑似ペアリングの根拠になった物体は画像を用いるべき
- ゲートの疑似教師情報



#### 定量的な実験結果

先行研究[Feng+, 2019][Laina+, 2019]と比較 → 提案手法はよりシンプルかつ高精度 BLEU-x, METEOR, ROUGE-L, CIDEr, SPICE...文生成の評価尺度 値が高い方が正解文に近い

BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr	SPICE
49.2	41.1	22.8	11.2	5.8	12.4	28.7	28.6
48.5	41.7	23.3	11.3	6.0	14.8	30.3	31.9
48.1	41.5	23.1	11.1	5.9	14.8	30.3	31.9

#### 定性的な実験結果

- ✓ ペアの教師情報がなくてもキャプション生成できる
- ✓ 先行研究[Feng+, 2019]と組合せると精度up

EACL 2021 採録!

## まとめ

### 複数の要求機能を満たすキャプション生成技術

- ✓ 行って帰ってくるロスによる詳細表現
- ✓ 疑似教師情報による教師データ不足への対応
- ✓ 材料統合ツリーを用いた作業動画のストーリー理解
- 今後の展望
  - ✓ コロナ禍で遅滞した個人の属性や傾向を含む総合データセットの整備 (クックパッド社との共同研究)
  - ✓ 個人への対応を含めた総合評価
  - ✓ 画像キャプション生成=異なるモダリティ間の変換 → 今後はより多様なモダリティ間のクロスモーダル変換へ

## ストーリーの多様性への対応

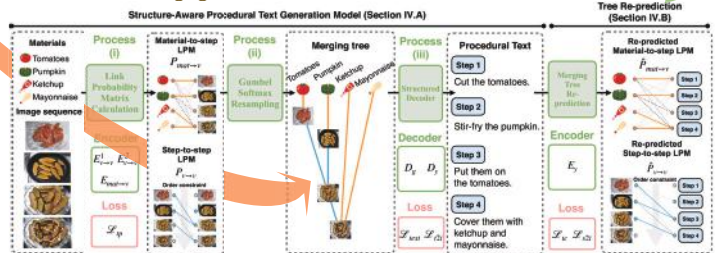
### 構造に配慮した写真列からの手順書生成

動画=主要な静止画の系列とみて  
材料+静止画系列 → レシピ生成

#### 画像列・レシピ、材料リストからなるデータセット vSIMMR の構築・公開

Datasets	Recipes?	Ingredients?	Structure?	Visual data	#Recipes
Breakfast [6]	✓	✓	✓	Video	N/A
EPIC-Kitchens [7]	✓	✓	✓	Video	N/A
YourCook2 [8]	✓	✓	✓	Image sequence	89
RecipeQA [9]	✓	✓	✓	Image sequence	19,799
Story boarding [1]	✓	✓	✓	Image sequence	16,405
Cookpad Image Dataset [10]	✓	✓	✓	Image sequence	1,715,595
Recipe Flow Graph [5]	✓	✓	✓	N/A	266
Action Graph [4]	✓	✓	✓	N/A	133
SIMMR [3]	✓	✓	✓	N/A	260
MM-Res [11]	✓	✓	✓	Image sequence	9,850
vSIMMR (ours)	✓	✓	✓	Image sequence	2,103

#### 材料の統合ツリー (Merging tree) からのレシピ生成+レシピからのツリー再推定



Ingredients	Step 1	Step 2	Step 3	Step 4	Step 5	Step 6
トマト	トマトを切る。	トマトを混ぜる。	トマトを炒める。	トマトを煮る。	トマトを焼く。	トマトを冷ます。
パンプキン	パンプキンを切る。	パンプキンを混ぜる。	パンプキンを炒める。	パンプキンを煮る。	パンプキンを焼く。	パンプキンを冷ます。
マヨネーズ	マヨネーズを切る。	マヨネーズを混ぜる。	マヨネーズを炒める。	マヨネーズを煮る。	マヨネーズを焼く。	マヨネーズを冷ます。

IEEE Access 採録!