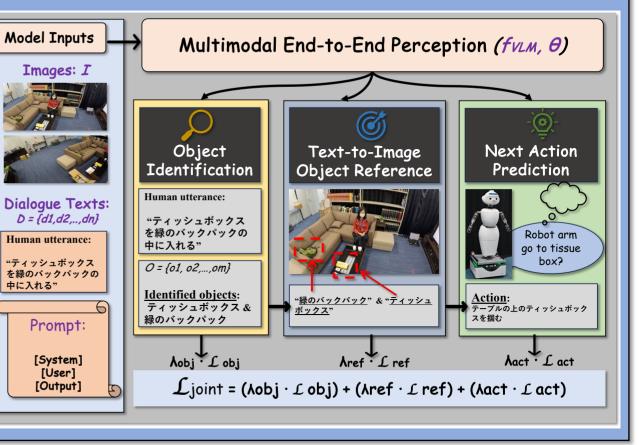


Attribute Annotation Template:

[category, color, shape, size, material, surface texture, position, state, function, brand/model, interactivity, proximity to the person]

```
"object": "Backpack".
                                                                                                                  "object": "Laptop".
"attributes": {
                                                                                                                  "attributes": {
   "Category": "個人アイテム".
                                                                                                                     "Category": "電子機器".
                                                                                                                     "Type": "ノートパソコン",
   "Color": "緑 (主要色)、黒 (副色)".
                                                                                                                     "Color": "黑",
   "Shape": "角が丸い長方形"。
                                                                                                                     "Shape": "長方形".
   "Size": "中型".
                                                                                                                     "Size": "中型".
   "Material": "布",
                                                                                                                     "Material": "プラスチック、金属".
   "Surface Texture": "滑らか",
                                                                                                                     "Surface Texture": "滑らか",
   "Position": "ソファの上",
                                                                                                                     "Position": "壁際の床の上".
   "State": "使用済み",
                                                                                                                     "State": "使用済み、画面オン".
   "Functionality": "収納、持ち運び",
                                                                                                                     "Functionality": "計算、情報表示",
   "Brand/Model": "指定なし",
   "Interactivity": "静的、可搬",
                                                                                                                     "Interactivity": "インタラクティブ、有線接続"
   "Proximity to Person": "近い、ソファの左アームの上"
                                                                                                                     "Proximity to Person": "遠い、壁際の床の上"
```



中に入れる"