

ALBEF

```
graph LR; A[ALBEF] --- B[相比ViLT，把image和text的特征提取器弄成了非对称的形式，image的特征提取器略大]; A --- C[在modality interaction之前利用CLIP的loss 训练了text和vision提取器的特征，使其更加一致]; A --- D[利用一个慢更新的动量模型提供pseudo-targets，解决了web data noisy的问题。就是网上的图片的文章描述可能不确切];
```

相比ViLT，把image和text的特征提取器弄成了非对称的形式，image的特征提取器略大

在modality interaction之前利用CLIP的loss 训练了text和vision提取器的特征，使其更加一致

利用一个慢更新的动量模型提供pseudo-targets，解决了web data noisy的问题。就是网上的图片的文章描述可能不确切