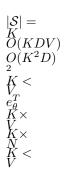
```
Con-
trastive
Weight
          ^{ing}_{2}/headless/imgs/hlm_{b}asic.pdfMaskedHeadlessLanguageModeling(HLM)usingContrastiveWeightTying.TheCWTotal National ContrastiveWeightTying.TheCWTotal National Contrastive National Contrastiv
            _{t}rain_{b}engio_{0}3, mnih 2012 fast, jean-\\
          etal-
2015-
          \begin{array}{c} \bar{u} \ddot{s} \ddot{i} \ddot{o} g, ma-\\ collins-\\ 2018-\\ T, \end{array}
               \widetilde{noise}. The semethod sapproximate the denominator of the soft max by using only a subset of the possible tokens. Those approach is a subset of the possible tokens. The semperature of the soft max by using only a subset of the possible tokens. The semperature of the soft max by using only a subset of the possible tokens. The semperature of the soft max by using only a subset of the possible tokens. The semperature of the soft max by using only a subset of the possible tokens. The semperature of the soft max by using only a subset of the possible tokens. The semperature of the soft max by using only a subset of the possible tokens. The semperature of the soft max by using only a subset of the possible tokens. The semperature of the soft max by using only a subset of the possible tokens. The semperature of the soft max by using the semperature of the semperature of
            Contrastive Estimation objective neethat use unique negative samples, contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach that samples representation of the contrary to our approach to our ap
            similarity to match pre-
            trained static embeddings for Machine Translation. We instead use the model's input embeddings a strainable target representation of the properties of the
            interspeech, wav2vec2, algayres-
       [merspeech, wav2vec2, digayres-etal-2022-dip. In NLP, contrastive learning has proven efficient in the training of sentence-level models gao-etal-2021-simcse, yan-etal-2021-consert, klein-nabi-2023-micse. Token-level approaches relyon contrastive auxiliary objectives that are added to the usual cross-entropyloss. SimCTGsu2022 contrastive introduces a token-level approaches relive bisections in single production of the pr
               level \c{c}ontrastive objective using in-
          batchoutput representations as negative samples, and adds this objective to a sentence—level contrastive loss and are gular causal LM loss. TaCLsu-et al-2022-tacl relies on a similar technique for encoder mod X = \begin{pmatrix} x_{i,j} \end{pmatrix}_{i \in [1,N], j \in [1,L]}
        \begin{array}{l} (x_{i,j})_{i \in [1,N], j \in [1,L]} \\ N \\ L \\ X \\ = \\ (\tilde{x}_{i,j})_{i \in [1,N], j \in [1,\tilde{L}]} \\ e_{\theta} \in \\ R^{V \times D} \\ V \\ T_{\theta} : \\ R^{N \times L \times D} \\ R^{N \times L \times D} \rightarrow \\ R^{P} \\ R^{P} \end{array} 
\begin{array}{l} R \in R \\ R \neq R \\ X_{\mathcal{S}} = \\ (x_{i,j})_{i,j \in \mathcal{S}} \\ (\tilde{x}_{i,j}) \\ T_{\theta} \\ X_{\mathcal{S}} \\ (\tilde{x}_{i,j}) \\ T_{\theta} \\ T_{\theta} \\ X_{\mathcal{S}} \\ (\tilde{x}_{i,j}) \\ T_{\theta} 
          \hat{p}_{i,j} = softmax \left( e_{\theta}^T \left( T_{\theta}(e_{\theta}(\tilde{x}_i))_j \right) \right)
       \mathcal{L}(\theta, X, \tilde{X}) = -\frac{1}{|\mathcal{S}|} \sum_{i,j \in \mathcal{S}} \mathbf{1}_{x_{i,j}} \cdot \log(\hat{p}_{i,j})
     \begin{array}{l} e_{\theta}^{T} \\ e_{\theta}(x_{i,j}) \\ o_{i,j}^{\theta} = \\ T_{\theta}(e_{\theta}(\tilde{x}_{i}))_{j} \\ e_{\theta}^{T} \\ \mathcal{S} \end{array}
          \mathcal{L}_{c}(\theta, X, \tilde{X}) = -\frac{1}{|\mathcal{S}|} \sum_{i,j \in \mathcal{S}} \frac{e^{o_{i,j}^{\theta} \cdot e_{\theta}(x_{i,j})}}{\sum_{k,l \in \mathcal{S}} e^{o_{i,j}^{\theta} \cdot e_{\theta}(x_{k,l})}}
                    \begin{array}{c} Con-\ trastive\ Weight \end{array}
               ing
per
se
do
```

 $\begin{array}{c} not \\ com-\\ bine \\ R^V \end{array}$ 



 $_2/headless/imgs/bert_train_speed_p13.pngTraininglatency \\_2/headless/imgs/bert_memory_use_p13.pngMemoryuse$ 

```
_{e}ncoder, we train a vanilla MLM and a headless counterpart. However, we share training hyperparameters such as batch size of the property of the property
 multienc.
Fine-tuned on English only
 48.0657.3274.0362.72 62 45.2552.1557.\overline{3}6\pm Fine-tuned on target language
\textbf{54.2566.9573.9669.1467.2260.0467.2265.5}^{\pm} \pm \\ curve_{m}ultilm.We find that the headless approach leads to significantly better performance for every language in both cross-line in the first f
  specific fine
  tuning. In average, the headless MLM outperforms its vanilla counterpart by 2 accuracy points in the cross-
 lingu\'{a}lscenario\', and by 2.7 points in the langu\'{a}ge-
  specific fine-
  tuning experiments.
 _2/headless/imgs/xnli_translate_train_p13.pngTranslate-Train: targetlanguagefine-tuning
 _2/headless/imgs/xnli_cross_p13.pngTranslate-Test:Englishfine-tuning
\begin{array}{l} -\frac{1}{c} urve_{m} ultilm, we evaluate the models at intermediate pretraining checkpoints and plotthe XNLI averages coreas a function of each property of the property o
{}_{t}rain_{m}ono_{e}nc. For each vocabulary size, we traina BPE to kenizer similar to the BERT to kenizer, and pretraina vanilla MLN {}_{2}/headless/imgs/bert_{v}ocab size_{p}13.png GLU E average score
 _2/headless/imgs/speed_vocabsize_p13.pngTrainingspeed\\
  _t oken count demonstrates that increasing the vocabulary size comes at almost node crease intraining speed for the HLMs, containing the vocabulary size comes at almost node crease in training speed for the HLMs, containing the vocabulary size comes at almost node crease in training speed for the HLMs, containing the vocabulary size comes at almost node crease in training speed for the HLMs, containing the vocabulary size comes at almost node crease in training speed for the HLMs, containing the vocabulary size comes at almost node crease in training speed for the vocabulary size comes at almost node crease in training speed for the vocabulary size comes at almost node crease in training speed for the vocabulary size comes at almost node crease in training speed for the vocabulary size comes at almost node crease in training speed for the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in the vocabulary size comes at almost node crease in t
  _{t}rain_{m}ono_{d}ec.HLMs are fine
  tunedsimilarlytosec: mono_decoder.
  _2/headless/imgs/batch_size_h^{\circ}ours_p 13.pngLAMBADA accuracy along pretraining for different batch sizes.
size, we observe that increasing batch size leads to better performance for our HLMs. While smaller batch size strain even faster ain mono_enc for different objectives. We observe that adding Cross-Entropy to CWT leads to slightly worse performance, at the cost of reduced throughput. We also notice that using a contrastive monotones are the cost of t
 x2.47
x2.1383.37
   "e
"But
what
  soft-
 max?"
head-
less
 mod-
el-
ing
 mostly
  pushes
  for
   crim-
    ħa-
  tion
be-
tween
  occurring
  to-kens
 ^{\circ}_{2}/headless/imgs/input_{e}mbs_{c}osines_{b}ert.pngMonolingualencoders
 _2/headless/imgs/input_embs_cosines_pythia.pngMonolingualdecoders
input_embs, we observe that HLM stendtogenerally represent synonym sinamore similar way than vanilla LMs, as cosine-similarity distributions slightly drift towards higher values. In average, cosine-similarity between synonym sis 1.4 point shigher for the encoder and roughly 7 point shigher for both the original HLM decoder and roughly 7 points higher for both the original HLM decoder and roughly 7 points higher for both the original HLM decoder and roughly 7 points higher for both the original HLM decoder and roughly 7 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the original HLM decoder and roughly 8 points higher for both the orig
 tunedversion.
bert-base-uncased
pythia-70m-deduped
 model-name
pythia-70m-deduped
pythia-70m-deduped
 (\beta_1,\beta_2) \\ \text{distilbert-base-multilingual-cased} \\ \text{distilbert-base-multilingual-cased}
distilbert-base-multilingual-cased google/bert_uncased L-4_H-512_A-8
google/bert_uncased_L-4_H-512_A-8 gpt2
         ize. The semodels rely on the GPT-
  \label{lem:action} \rat{2} architecture with a few changes. The sechanges scaledown the model size to 11 M parameters.
```

 $(w_c)_{c \in [1,C]}$ 

 $\mathcal{L}_{ce}(X,Y) = \sum_{c=1}^{C} \mathcal{L}_{c}(X,Y)$   $\mathcal{L}_{c}(X,Y) = \sum_{i=1}^{N} \mathbf{1}_{y_{i}=c} \cdot \mathcal{L}_{ce}(x_{i},y_{i})$