

# Backpack Language Models

John Hewitt John Thickstun Christopher D. Manning Percy Liang

Department of Computer Science, Stanford University  
{johnhew, jthickstun, manning, plieng}@cs.stanford.edu

## Abstract

We present *Backpacks*: a new neural architecture that marries strong modeling performance with an interface for interpretability and control. Backpacks learn multiple non-contextual *sense* vectors for each word in a vocabulary, and represent a word in a sequence as a context-dependent, non-negative linear combination of sense vectors in this sequence. We find that, after training, sense vectors specialize, each encoding a different aspect of a word. We can interpret a sense vector by inspecting its (non-contextual, linear) projection onto the output space, and intervene on these interpretable hooks to change the model’s behavior in predictable ways. We train a 170M-parameter Backpack language model on OpenWebText, matching the loss of a GPT-2 small (124M-parameter) Transformer. On lexical similarity evaluations, we find that Backpack sense vectors outperform even a 6B-parameter Transformer LM’s word embeddings. Finally, we present simple algorithms that intervene on sense vectors to perform controllable text generation and debiasing. For example, we can edit the sense vocabulary to tend more towards a topic, or localize a source of gender bias to a sense vector and globally suppress that sense.

## 1 Introduction

Consider the prefix *The CEO believes that* \_\_\_\_, and the problem of debiasing a neural language model’s distribution over *he/she*. Intuitively, the bias for *he* originates in the word *CEO*, because replacing *CEO* with *nurse* flips the observed bias. A successful intervention to debias *CEO* must reliably apply in all contexts in which the word *CEO* appears; ideally we would want to make a **non-contextual** change to the model that has predictable effects in **all contexts**. In general, in all aspects of interpretability and control, it is desirable to make interventions with a tractable interface (e.g., non-contextual representations) that apply globally.

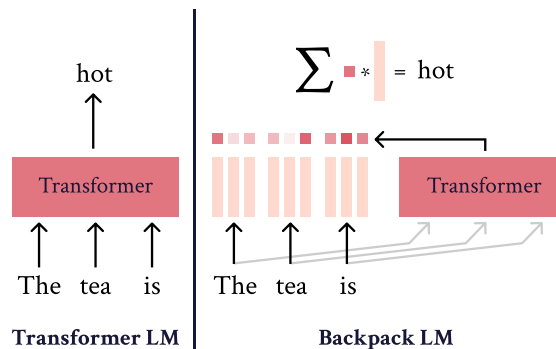


Figure 1: Transformers are monolithic functions of sequences. In Backpacks, the output is a weighted sum of non-contextual, learned word aspects.

Such interventions are difficult in Transformer models (Vaswani et al., 2017) because their contextual representations are monolithic functions of their input. Almost any intervention on the model has complex, non-linear effects that depend on context. We would instead like models that enable precise, rich interventions that apply predictably in all contexts, and are still expressive, so they are a viable alternative to Transformers.

We address these challenges with a new neural architecture, the *Backpack*, for which predictions are log-linear combinations of non-contextual representations. We represent each word in a vocabulary as a set of non-contextual *sense vectors* that represent distinct learned aspects of the word. For example, sense vectors for the word “science” could encode types of science, connections to technology, notions of science being “settled,” or different aspects of the scientific process (replication or experiment) (Table 1). Sense vectors do not learn classic word sense, but more general aspects of a word’s potential roles in different contexts; in fact, they can be seen as a multi-vector generalization of classic word vectors (Mikolov et al., 2013).<sup>1</sup>

<sup>1</sup>Our code, sense vectors, language model weights, and demos are available at <https://backpackmodels.science>.

A few senses of the word <i>science</i>				
Sense 3	Sense 7	Sense 9	Sense 10	Sense 8
fiction	replication	religion	settled	clones
fictional	citation	rology	sett	experiments
Fiction	Hubble	hydra	settle	mage
literacy	reprodu	religions	unsett	experiment
denial	Discovery	nec	Sett	rats

$$\mathit{MacBook}_{HP} = \mathit{MacBook} - \mathit{Apple} + \mathit{HP}$$

**The MacBook is best known for** its form factor, but HP has continued with its Linux-based computing strategy. HP introduced the Hyper 212 in 2014 and has continued to push soon-to-be-released 32-inch machines with Intel’s Skylake processors.

Table 1: Examples of the rich specialization of sense vectors representing the word *science*, and an example of editing sense vectors non-contextually (changing MacBook to be associated with HP) and having the resulting *contextual* predictions change.

To make interventions on sense vectors behave predictably in different contexts, a Backpack represents each word in a sequence as a **linear combination** of the sense vectors for all words in the sequence. The expressivity of a Backpack comes from the network that computes the weights of the linear combination as a function of the whole sequence; for example, in all our experiments we use a Transformer for this. Since sense vectors are softly selected depending on the context, they can specialize; each sense can learn to be predictively useful in only some contexts. The log-linear contribution of senses to predictions then implies that the interventions on sense vectors we demonstrate in Section 6 apply identically (up to a non-negative scalar weight) regardless of context.

Our experiments demonstrate the expressivity of Backpack language models, and the promise of interventions on sense vectors for interpretability and control. In Section 4 we train Backpack language models on 50B tokens (5 epochs) of OpenWebText; a Backpack with 124M parameters in the contextual network (and 46M parameters for sense vectors) achieves the perplexity of a 124M-parameter Transformer; thus one pays for more interpretability with a larger model size. In Section 5, we show that sense vectors specialize to encode rich notions of word meaning. Quantitatively, on four lexical similarity datasets (e.g., SimLex999), sense vectors of a 170M parameter Backpack outperform word embeddings of the 6B-parameter GPT-J-6B Transformer, and approach the performance of state-of-the-art specialized methods for this task. Finally, in Section 6 we show that sense vectors offer a control mechanism for Backpack language models. For example, stereotypically gendered profession words (e.g., “CEO” or “nurse”) tend to learn a sense vector associated with this gender bias; by downscaling this sense vector, we greatly reduce disparity in contextual predictions in a limited setting.

## 2 The Backpack Architecture

In this section, we define the general form of the Backpack architecture. We then show how continuous bag-of-words word2vec (CBOV) (Mikolov et al., 2013) and Self-Attention-Only networks (Elhage et al., 2021; Olsson et al., 2022) are special cases of Backpacks.

### 2.1 Backpack General Form

A Backpack is a parametric function that maps a sequence of symbols  $\mathbf{x}_{1:n} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  to a sequence of vectors  $\mathbf{o}_{1:n} = (\mathbf{o}_1, \dots, \mathbf{o}_n)$ , where each symbol  $\mathbf{x}_i$  belongs to a finite vocabulary  $\mathcal{V}$  and  $\mathbf{o}_i \in \mathbb{R}^d$ . We call  $\mathbf{o}_i$  the *Backpack representation* of  $\mathbf{x}_i$  in the context of a sequence  $\mathbf{x}_{1:n}$ .

**Sense vectors.** For each  $\mathbf{x} \in \mathcal{V}$ , a Backpack constructs  $k$  sense vectors

$$C(\mathbf{x})_1, \dots, C(\mathbf{x})_k, \quad (1)$$

where  $C : \mathcal{V} \rightarrow \mathbb{R}^{k \times d}$ . Sense vectors are a multi-vector analog to classic non-contextual word representations like word2vec or GloVe: we make this analogy precise in Section 2.2.

**Weighted sum.** For a sequence  $\mathbf{x}_{1:n}$ , the representation  $\mathbf{o}_i$  of element  $\mathbf{x}_i$  is a weighted sum of the predictive sense vectors for the words in its context: given *contextualization weights*  $\alpha \in \mathbb{R}^{k \times n \times n}$ ,

$$\mathbf{o}_i = \sum_{j=1}^n \sum_{\ell=1}^k \alpha_{\ell ij} C(\mathbf{x}_j)_\ell. \quad (2)$$

The contextualization weights  $\alpha_{\ell ij}$  of a Backpack are themselves defined by a (non-linear) *contextualization function* of the entire sequence  $\mathbf{x}_{1:n}$ :

$$\alpha = A(\mathbf{x}_{1:n}), \quad (3)$$

where  $A : \mathcal{V}^n \rightarrow \mathbb{R}^{k \times n \times n}$ .

The name ‘‘Backpack’’ is inspired by the fact that a backpack is like a bag—but more orderly. Like a bag-of-words, a Backpack representation is a sum of non-contextual senses; but a Backpack is more orderly, because the weights in this sum depend on the ordered sequence.

**Backpack Models.** A *Backpack model* is a probabilistic model that defines probabilities over some output space  $\mathcal{Y}$  as a log-linear function of a Backpack representation  $\mathbf{o}_{1:n} \in \mathbb{R}^{n \times d}$ :

$$p(\mathbf{y} | \mathbf{o}_{1:n}) = \text{softmax}(E(\mathbf{o}_{1:n})), \quad (4)$$

where  $\mathbf{y} \in \mathcal{Y}$  and  $E : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{|\mathcal{Y}|}$  is a linear transformation. Because Backpack models are log-linear in their representations, the sense vectors contribute log-linearly to predictions. This allows us to inspect a sense vector by projecting it onto the vocabulary via  $E$  and observe exactly how it will contribute to predictions in any context.

Models parameterized by the prevailing deep neural architectures—including LSTMs (Hochreiter and Schmidhuber, 1997) and Transformers—are not Backpacks because their output representations are (relatively) unconstrained functions of the entire sequence. By contrast, Backpack models may seem limited in expressivity: the representations  $\mathbf{o}_j$  are scalar-weighted sums of non-contextual vectors  $C(\mathbf{x}_j)_\ell$ . Contextual relationships between sequence elements can only be expressed through the weights  $\alpha = A(\mathbf{x}_{1:n})$ . Nevertheless, our experiments show that an expressive contextualization weight network can represent complex functions by weighted sums of sense vectors, e.g., our 170M parameter Backpack LM uses a 124M-parameter Transformer to compute  $\alpha$ , and achieves the loss of a 124M-parameter Transformer LM.

To place Backpacks in some historical context, we now show how two existing architectures can be described as Backpacks.

## 2.2 Continuous Bag-of-Words is a Backpack

The continuous bag-of-words word2vec model defines a probability distribution over a center word  $\mathbf{x}_c \in \mathcal{V}$  conditioned on  $n$  context words  $\mathbf{x}_{1:n}$ .<sup>2</sup> The model proceeds to (1) construct vector embeddings  $\mathbf{v}_\mathbf{x}$  for each  $\mathbf{x} \in \mathcal{V}$ , and (2) uniformly average the embeddings of the context words to predict the

<sup>2</sup>Context in this setting is usually defined as words surrounding the center word.

center word:

$$\bar{\mathbf{v}}_{\mathbf{x}_c} = \sum_{i=1}^n \frac{1}{n} \mathbf{v}_{\mathbf{x}_i}, \quad (5)$$

$$p(\mathbf{x}_c | \mathbf{x}_{1:n}) = \text{softmax}(U\bar{\mathbf{v}}_{\mathbf{x}_c}), \quad (6)$$

where  $U \in \mathbb{R}^{\mathcal{V} \times d}$ . We see that  $\bar{\mathbf{v}}_{\mathbf{x}_c}$  is a Backpack representation by setting  $C(\mathbf{x}) = \mathbf{v}_\mathbf{x} \in \mathbb{R}^{1 \times d}$  in Equation (1) using a single sense vector ( $k = 1$ ) and setting the contextualization weights in Equation (3) to be uniform:  $\alpha_{\ell ij} = \frac{1}{n}$ .

This connection to CBoW foreshadows the emergence of linguistic structures in the predictive sense vectors of Backpack models, just as these structures emerge in CBoW (Mikolov et al., 2013).

## 2.3 Single-Layer Self-Attention is a Backpack

The Backpack structure—define sense vectors (values), and use the sequence to determine how to sum them (weights)—may remind the reader of a single layer of self-attention. The key-query-value self-attention function is as follows:

$$\mathbf{o}_j = \sum_{i=1}^n \sum_{\ell=1}^k \alpha_{\ell ij} O V^{(\ell)} \mathbf{x}_j \quad (7)$$

$$\alpha_\ell = \text{softmax}(\mathbf{x}^\top K^{(\ell)\top} Q^{(\ell)} \mathbf{x}), \quad (8)$$

where  $\mathbf{x} \in \mathbb{R}^{n \times d}$  is (overloaded) to be a non-contextual embedding of the sequence,  $O \in \mathbb{R}^{d \times d/k}$ , and  $V^{(\ell)} \in \mathbb{R}^{d/k \times d}$ , where  $k$  is the number of attention heads. The self-attention function is a Backpack with  $C(\mathbf{x}_j)_\ell = O V^{(\ell)} \mathbf{x}_j$ . Self-attention-only networks are studied in the context of, e.g., mechanistic interpretability (Elhage et al., 2021). A Transformer composes blocks of self-attention and non-linear feed-forward layers that combine information from the whole sequence; unlike a Transformer, the contextualization weights of a Backpack each select a non-contextual sense of a single word.

## 3 Language Modeling with Backpacks

In this section, we define a neural autoregressive language model parameterized by a Backpack. We use the standard softmax parameterization of the probability over the next token in a sequence, with a weight matrix  $E \in \mathbb{R}^{d \times |\mathcal{V}|}$  that maps a representation  $\mathbf{o}_j \in \mathbb{R}^d$  to logits  $E^\top \mathbf{o}_j \in \mathbb{R}^{|\mathcal{V}|}$ :

$$p(\mathbf{x}_j | \mathbf{x}_{1:j-1}) = \text{softmax}(E^\top \mathbf{o}_j). \quad (9)$$

Recall (Section 2.1) that Backpack representations  $\mathbf{o}_j$  are defined by sense vectors  $C(\mathbf{x})$  and contextualization weights  $\alpha_j$ . In Section 3.1 we describe a parameterization of  $C$  for the predictive sense vectors in Equation (1), and in Section 3.2 we describe a parameterization of  $A$  for the contextualization weight network in Equation (3). When  $\mathbf{o}_j$  is parameterized by a Backpack, we call a model of the form given by Equation (9) a *Backpack LM*.

### 3.1 Parameterizing senses

For the sense function  $C : \mathcal{V} \rightarrow \mathbb{R}^{k \times d}$ , we embed each  $\mathbf{x} \in \mathcal{V}$  into  $\mathbb{R}^d$  and pass these embeddings through a feed-forward network  $\text{FF} : \mathbb{R}^d \rightarrow \mathbb{R}^{k \times d}$ :

$$C(\mathbf{x}) = \text{FF}(E\mathbf{x}), \quad (10)$$

where the embedding/projection matrix  $E$  is tied to the output matrix in Equation (9) (Press and Wolf, 2017). Note that we could define all  $k \times |\mathcal{V}|$  sense vectors using a lookup table, but this would be an enormous number of parameters as  $k$  grows large. Instead, we embed the words as  $E\mathbf{x} \in \mathbb{R}^d$ , and then blow them up to  $\mathbb{R}^{d \times k}$  using shared weights. This may explain the related sense roles observed for different word types in Section 5.1.

### 3.2 Parameterizing contextualization weights

We parameterize  $A : \mathcal{V}^n \rightarrow \mathbb{R}^{k \times n \times n}$  using a standard Transformer, followed by a layer of multi-headed key-query self-attention. That is, we pass an embedded sequence through a Transformer

$$\mathbf{h}_{1:n} = \text{Transformer}(E\mathbf{x}_{1:n}) \quad (11)$$

(with proper autoregressive masking and some position representation) and compute  $A(\mathbf{x}_{1:n}) = \alpha$ , where

$$\alpha_\ell = \text{softmax}(\mathbf{h}_{1:n} K^{(\ell)\top} Q^{(\ell)} \mathbf{h}_{1:n}^\top), \quad (12)$$

for each predictive sense  $\ell = 1, \dots, k$  with matrices  $K^{(\ell)}, Q^{(\ell)} \in \mathbb{R}^{d \times d/k}$ . We can think of the  $k$  senses as heads and, for each head, the contextualization weights define a distribution of attention over words.<sup>3</sup>

## 4 Experiments Training Backpack LMs

In this section we specify the hyperparameters used to train Backpack and Transformer language models (Section 4.1), data and optimization procedure

<sup>3</sup>Note that the sense weights are normalized (1) independently for each sense, and (2) to sum to one over the sequence length.

(Section 4.2), evaluations (Section 4.3) and results (Section 4.4). We also show the necessity of learning  $k > 1$  sense vectors to achieve strong language modeling performance (Section 4.5).

### 4.1 Models

We train three Transformer baseline models, which we label Micro (30M parameters), Mini (70M parameters), and Small (124M parameters; the same size as GPT-2 small). We also train Micro (40M), Mini (100M), and Small (170M) Backpack language models, for which the weighting function (Equation 11) is parameterized using the corresponding Transformer, and almost all extra parameters are in the non-contextual sense vectors.<sup>4</sup> Backpacks thus cost extra parameters and compute beyond their underlying contextualization network. Except where stated, we use  $k = 16$  sense vectors in all Backpacks (Section A).

We use a reduced sequence length of 512 for all models, and the 50,257-subword GPT-2 tokenizer. Model hidden dimensionalities, layer counts, and head counts are reported in Table 9.

### 4.2 Data & Optimization

We train all models on OpenWebText (Gokaslan and Cohen, 2019), a publicly available approximate reconstruction of the English WebText corpus used to train the GPT-2 family of models (Radford et al., 2019). We use a batch size of 524,288 tokens, and train all models for 100,000 gradient steps for a total of 52B tokens; training for longer is known to make marginal difference for small models (Hoffmann et al., 2022). The size of OpenWebText means this is roughly 5 epochs. We use cross-entropy loss and the AdamW optimizer, with a warmup of 5,000 steps and linear decay to zero.

### 4.3 Evaluations

Before our experiments in interpretability and control, we check the expressivity of Backpacks. We evaluate models on perplexity for a held out set of OpenWebText, perplexity and accuracy for the (OpenAI variant of) LAMBADA evaluation of long-distance dependencies (Radford et al., 2019; Paperno et al., 2016), perplexity on Wikitext (Merity et al., 2017), and BLiMP English linguistic competence accuracy (Warstadt et al., 2020) evaluated using the EleutherAI harness (Gao et al., 2021) (Version 1).

<sup>4</sup>There are a negligible number of additional parameters in the final key-query Backpack operation (Equation 12)).

Model	OpenWebText PPL ↓	LAMBADA PPL ↓	LAMBADA ACC ↑	Wikitext PPL ↓	BLiMP ↑
Backpack-Micro	<b>31.5</b>	<b>110</b>	<b>24.7</b>	<b>71.5</b>	75.6
Transformer-Micro	34.4	201	21.3	79.5	<b>77.8</b>
Backpack-Mini	<b>23.5</b>	<b>42.7</b>	<b>31.6</b>	<b>49.0</b>	76.2
Transformer-Mini	24.5	58.8	29.7	52.8	<b>80.4</b>
Backpack-Small	<b>20.1</b>	<b>26.5</b>	<b>37.5</b>	<b>40.9</b>	76.3
Transformer-Small	<b>20.2</b>	32.7	34.9	42.2	<b>81.9</b>

Table 2: Language modeling performance; all models trained for 100k steps, 500K token batch size, on OWT. For PPL, lower is better; for accuracy, higher is better. Note that models are not parameter-comparable; each Backpack has a matched-size Transformer in its contextualization network.

#### 4.4 Discussion

Comparing each Backpack LM to a Transformer LM of equivalent specification to the Backpack’s contextualization network, we see that the Backpack performs roughly as well (Table 2). Again, the Backpack has more parameters, a tax for the interface provided by sense vectors. During training, we find that Backpack language models take longer to converge than Transformers. Curiously, while the Small Backpack and Transformer achieve almost identical OWT perplexity, the Backpack language models perform substantially better on LAMBADA and Wikitext, but worse on BLiMP.

#### 4.5 Effect of varying the number of senses

To study the impact of the number of sense vectors on language modeling performance, we train Mini-sized Backpack language models on a reduced schedule of 50,000 gradient steps, for  $k \in \{1, 4, 16, 64\}$  sense vectors. The perplexities for  $k = 1, 4, 16, 64$  are 38.6, 29.3, 26.0, and 24.1, demonstrating the necessity of a non-singleton set of sense vectors. Table 8 contains the full results.

### 5 Emergent Structure in Sense Vectors

Backpack language model sense vectors are not trained using a supervised notion of word sense, but implicitly specialize to encode different shades of a word’s predictive use. In this section, we qualitatively examine sense vectors (Section 5.1) and quantitatively demonstrate their effectiveness in computing lexical similarity and relatedness (Section 5.2). Taken together, this suggests that sense vectors can provide a high-level interface for intervention, which we explore in Section 6.

#### 5.1 Visualizing Senses

Empirically, trained Backpack models associate specific sense vector indices with different roles for

prediction. We interpret these roles by picking a sense  $\ell$  of a word  $\mathbf{x}$ , and projecting this sense onto the word embeddings:  $E^T C(\mathbf{x})_\ell \in \mathbb{R}^{|\mathcal{V}|}$ . Note that this is *exactly* (up to a scalar) how this sense contributes to any prediction of the model. We interpret a sense vector’s role by reporting the words with the highest score under this projection.

Table 3 visualizes a few of these senses. For example, sense 12 seems to encode a broad notion of relatedness for almost all words; sense 3 encodes particulars of the bigram distribution given  $\mathbf{x}$ ; sense 14 seems to encode both associated objects for verbs, and noun modifier dependency children for nouns. In Section 5.2 we show that sense 14 encodes a powerful notion of verb similarity.

#### 5.2 Lexical Relationship Tests

Classic lexical-relatedness and similarity tests measure the extent to which a similarity function on pairs of words correlates with human-elicited notions of similarity. Similarity functions derived from word embeddings are evaluated by Spearman correlation between the predicted and true similarity rank-order. Early non-contextual embeddings like COALS (Rohde et al., 2005), word2vec (Mikolov et al., 2013), and GloVe (Pennington et al., 2014) have recently been outperformed by word embeddings derived by distillation of contextual networks (Bommasani et al., 2020; Gupta and Jaggi, 2021; Chronis and Erk, 2020). We evaluate Backpack LM sense vectors on similarity datasets SimLex999 (Hill et al., 2015), SimVerb3500 (Gerz et al., 2016), and relatedness datasets RG65 (Rubenstein and Goode-nough, 1965) and (Agirre et al., 2009).

**Sense $_\ell$  Cosine.** For all  $\ell \in \{1, \dots, k\}$ , we define a similarity function based only on sense  $\ell$ :

$$\text{Sim}_\ell(\mathbf{x}, \mathbf{x}') = \text{cossim}(C(\mathbf{x})_\ell, C(\mathbf{x}')_\ell), \quad (13)$$

Sense 12 ( <i>relatedness</i> )				Sense 14 ( <i>Verb objects, nmod nouns</i> )			
<i>tasty</i>	<i>quickly</i>	<i>Apple</i>	<i>believe</i>	<i>build</i>	<i>attest</i>	<i>importance</i>	<i>appreciate</i>
tasty	quick	Apple	belief	bridges	worthiness	maintaining	finer
culinary	quickest	Apple	Belief	wall	Published	wellbeing	nuance
tasted	quick	iPhone	beliefs	lasting	superiority	teamwork	beauty
delicious	quicker	iPhone	believing	ig	accuracy	plurality	irony
taste	fast	iPhones	believe	rapport	validity	upholding	simplicity

Sense 3 ( <i>next wordpiece</i> )			Sense 7 ( <i>Proper Noun Associations</i> )		
<i>pizza</i>	<i>interest</i>	<i>the</i>	<i>Apple</i>	<i>Obama</i>	<i>Messi</i>
cutter	rate	slightest	macOS	Dreams	Messi
tracker	rates	same	iCloud	Barack	Argentina
iol	groups	entirety	Siri	Ob	Mess
makers	waivers	rest	iOS	Michelle	Barcelona
maker	waiver	latter	tv	Jeremiah	iesta

Table 3: Visualization of how the same sense index across many words encodes fine-grained notions of meaning, relatedness, and predictive utility. Each sense is given a label thought up by the authors, and for a few words, the target words that are highest scored by the sense vector.

Model	SL999	SV3500	RG65	WS353
<i>Classic Non-Contextual Embeddings</i>				
word2vec	0.442	0.367	0.679	0.684
GloVe	0.371	0.227	0.687	0.607
<i>Embeddings from large existing models</i>				
GPT2-1.5B	0.523	0.418	0.670	0.706
GPT-J-6B	0.492	0.374	<b>0.766</b>	0.673
<i>Embeddings from our models + baseline Transformer</i>				
Trnsf 124M	0.478	0.363	0.634	0.681
Sim <sub>12</sub> (ours)	0.522	0.471	0.754	<b>0.749</b>
Sim <sub>14</sub> (ours)	0.500	<b>0.502</b>	0.591	0.655
Sim <sub>min</sub> (ours)	<b>0.540</b>	0.471	0.653	0.607
<i>Special-purpose SOTA models</i>				
SOTA (Single)	0.554	0.473	0.835	0.764
SOTA (Multi)	0.605	0.528	-	0.807

Table 4: Results on lexical similarity evaluation. All numbers are Spearman correlations; higher is better.

where  $\text{cossim}$  is cosine similarity. Intuitively, we expect that some senses may specialize to learn lexical relatedness or similarity.

**Minimum Sense Cosine.** Because each sense encodes a different aspect of a word’s meaning, we might expect that highly similar words are similar across *all* senses. We test for this strong form of similarity using

$$\text{Sim}_{\min}(\mathbf{x}, \mathbf{x}') = \min_{\ell} \text{Sim}_{\ell}(\mathbf{x}, \mathbf{x}') \quad (14)$$

**Other methods.** We evaluate embeddings from the tied softmax/embedding matrices of the much larger GPT-2-1.5B (Radford et al., 2019) and GPT-J-6B (Wang and Komatsuzaki, 2021), classic word embeddings (from Bommasani et al. (2020)) and

state-of-the art specialized methods using either a single vector per word (Gupta, 2021) or many vectors (Chronis and Erk, 2020).

**Discussion.** Sense 12 (the “synonym” sense) performs well across datasets, matching or outperforming embeddings like GPT-2-1.5B and GPT-J-6B (Except GPT-J-6B on RG-65). Sense 14, the “verb objects” sense, performs best on just verb similarity (VerbSim3500), and the minimum similarity over senses works especially well on noun lexical similarity (SimLex999.) Our methods approach the performance of state-of-the-art methods; despite being trained for a very different task, sense vectors encode substantial lexical information (Table 4).

## 6 Sense Vectors for Control

In this section, we demonstrate several proof-of-concept methods that leverage sense vectors for controlling LM behavior.

### 6.1 Topic-controlled generation

Given a bag-of-words target  $b \in \mathbb{R}^{|\mathcal{V}|}$ , e.g., *arts*, *culture*, we would like to bias generation towards sequences related to concepts related to these terms. Our algorithm proceeds in three parts. First, we sort sense vectors by log-probability assigned to  $b$ , that is,  $b^{\top}(E^{\top}C(\mathbf{x})_{\ell})$ .<sup>5</sup> Second, based on the scores, we assign a re-weighting factor  $\delta$  to each sense; senses with the higher scores weighted more. (See Section D for details.) Third, we generate from

<sup>5</sup>We divide this term by the maximum absolute log-probability of the sense vector,  $\max_{x \in \mathcal{V}} \mathbf{x}^{\top}(E^{\top}C(\mathbf{x})_{\ell})$ .

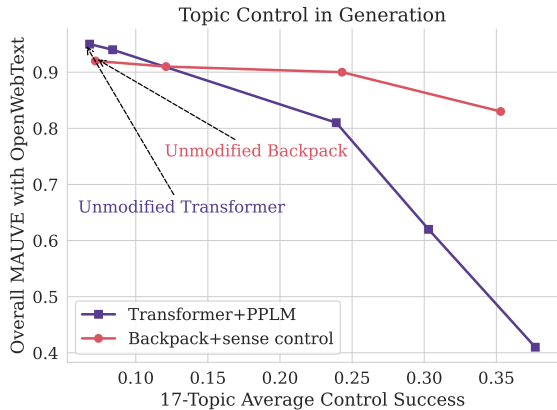


Figure 2: Results in controlling topic via sense intervention in Backpacks, and PPLM in Transformers.

the Backpack using the re-weighted sense vectors, reducing  $\delta$  back to 1 as the topic is introduced. The updated backpack equation is

$$\mathbf{o}_i = \sum_{j=1}^n \sum_{\ell=1}^k \alpha_{\ell ij} \delta_{\ell ij} C(\mathbf{x}_j)_\ell, \quad (15)$$

where  $\delta_{ij\ell}$  is the re-weighting. Intuitively, the semantic coherence of sense vectors may imply that upweighting senses with affinity to the target bag-of-words richly upweights related words and topics. We give details as to how we perform the sense re-weighting and the annealing in Section D.

**Evaluation.** We use the label descriptors of the topic classifier of Antypas et al. (2022), with 17 categories (*sports, arts & culture, health,...*), as the bag-of-words for control. We evaluate control accuracy as the percent of generations to which the classifier assigns the correct topic label, and overall generation quality and diversity using MAUVE scores (Pillutla et al., 2021).<sup>6</sup>

**Results.** We compare to Plug-and-Play Language Models (PPLM; Dathathri et al. (2019)), a considerably slower, gradient-based control method using our Small Transformer model. We generate 500 samples from each model for each topic across a range of strengths of control. We find that sense controlled generation provides at least as strong control as PPLM (Figure 2), though the MAUVE scores of the unmodified Transformer are higher than the Backpack.) Results and examples are provided in the Appendix in Tables 12, 16, 17, 18.

<sup>6</sup>We concatenate generations across the 17 categories and compute MAUVE against OpenWebText validation examples.

Model	Bias Ratio ↓	Reduction %
Unbiased	1	-
<i>Transformer</i>		
Unmodified	7.02	-
Project-Nullspace	6.72	5%
Optimize-Nullspace	7.02	0%
<i>Backpack</i>		
Unmodified	4.34	-
Remove-Sense10	2.88	44%
Optimize-Sense10	2.16	65%

Table 5: Pronoun-based gender bias reduction in a limited setting.

## 6.2 Mitigating gender bias

Through inspection, we learned that sense vector 10 of many stereotypically gendered profession nouns (nurse, CEO, teacher) coherently express the stereotype through pronouns. Table 13 gives examples of these senses. We attempt to mitigate gender bias in Backpack behavior on these gendered profession nouns by *turning down* sense 10 (multiplying by a scalar less than 1).

We took an existing set of stereotypically gendered profession nouns from WinoBias (Zhao et al., 2018), and constructed a simplified setting in which a single profession word is in each context, and a third-person nominative pronoun (e.g., he/she/they) is acceptable, e.g., *My CEO said that\_\_*. The full set of nouns and prompts is in Section D.2. We evaluate models on the average of the bias of probabilities of *him* vs *her* as follows:

$$\mathbb{E}_{\mathbf{x} \in \text{prompts}} \left[ \max \left( \frac{p(\text{he} | \mathbf{x})}{p(\text{she} | \mathbf{x})}, \frac{p(\text{she} | \mathbf{x})}{p(\text{he} | \mathbf{x})} \right) \right].$$

**Baseline.** To debias a Transformer with an analogous method, we take inspiration from Bolukbasi et al. (2016). We take  $Ex_{\text{he}} - Ex_{\text{she}}$  as an estimate of a gender bias direction, and project the embedding  $Ex_{\text{nurse}}$  either to the nullspace of this direction or only partially remove it.

**Results.** A perfectly unbiased model would achieve ratio 1, whereas the unmodified Transformer achieves 7, and with nullspace projection, 6.72 (Table 5). Finding the optimal fraction of the gender bias direction to remove per profession does not improve further. For Backpacks, we find that removing sense 10 from the profession word (setting it to zero) reduces the bias score from 4.34 to 2.88. Learning the optimal removal fraction per profession achieves 2.16, for a total reduction of

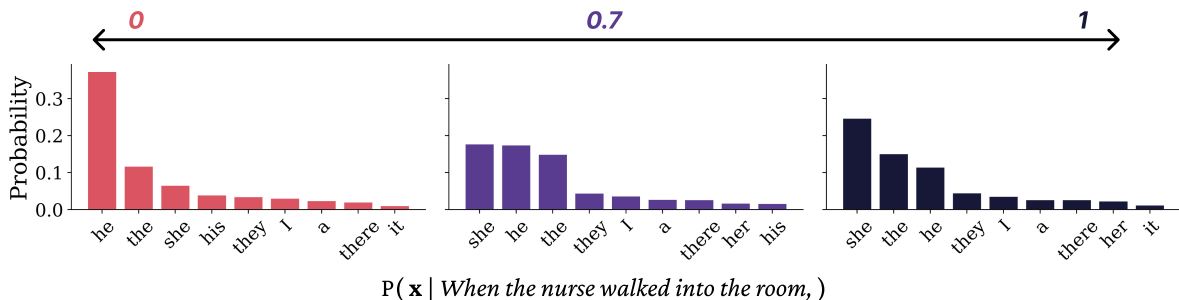


Figure 3: The effect on the conditional probability distribution of a Backpack LM on the prefix *when the nurse walked into the room,* of modulating the effect of sense 10 of *nurse* from 0 (totally removed) to 1 (original.)

---

**The MacBook is best known for** its form factor, but HP has continued with its Linux-based computing strategy. HP introduced the Hyper 212 in 2014 and has continued to push soon-to-be-released 32-inch machines with Intel’s Skylake processors.

---

**The MacBook didn’t come into the picture until 2000,** when HP followed up with a 15-year flood of HP available laptops.

---

**I was thinking about Brady’s role on** the Colts before joining other high-profile signings. This is what McElhaney and I discussed.

McElhaney: Look, what I didn’t mean by this is we didn’t move. We think that we’re getting a lot better, too.

---

Table 6: Samples from a Backpack wherein *Apple* has been projected out of the *MacBook* sense embeddings, and replaced with *HP*. Likewise with *Brady*, *Patriots*, and *Colts*. Prompts are bolded.

65%.<sup>7</sup> In Figure 3, we demonstrate the clear effect of ablating sense 10 on the most likely words in one of these contexts.<sup>8</sup>

### 6.3 Knowledge editing

Sense vectors show promise for use in *knowledge editing* (De Cao et al., 2021)—editing a model’s predictions about world knowledge. In particular, many associations with proper nouns can be localized to sense vectors in that noun. In this qualitative proof-of-concept, we edit the sense vectors of a target word  $\mathbf{x}$  (e.g., *MacBook* to remove associations with a word  $\mathbf{x}_r$  (e.g., *Apple*) and replace those associations with another word  $\mathbf{x}_a$  (e.g., *HP*). Intuitively, this intervention ensures that whenever the contextualization weights would point to a sense vector in *MacBook* to predict words associated with *Apple*, it now predicts words associated with *HP*.

<sup>7</sup>Curiously, Backpacks are overall less biased to begin with (in this setting); we don’t have a strong hypothesis as to why.

<sup>8</sup>It is incidental that sense 10 encodes gender bias as opposed to another sense index; the consistency in index across words may be due to parameter sharing in  $C$ .

We project each sense vector of  $\mathbf{x}$  to the nullspace of  $E\mathbf{x}_r$ , and then add in  $E\mathbf{x}_a$ :

$$\tilde{C}(\mathbf{x})_\ell = C(\mathbf{x})_\ell + \frac{C(\mathbf{x})_\ell^\top E\mathbf{x}_r}{\|C(\mathbf{x}_r)_\ell\|_2^2} \left( \frac{E\mathbf{x}_a}{\phi} - E\mathbf{x}_r \right),$$

where  $\phi = \frac{\|E\mathbf{x}_a\|_2^2}{\|E\mathbf{x}_r\|_2^2}$  is a normalization term to account for the differing norms of  $E\mathbf{x}_a$  and  $E\mathbf{x}_r$ . Intuitively, this projection modifies each sense vector in measure proportional to how much  $\mathbf{x}_r$  was predicted by that sense. So, senses of *MacBook* that would added mass to *Apple* now add mass to *HP*; unrelated senses are not affected. In Table 6, we show samples providing intuition for how *MacBook* evokes HP instead of Apple, but is otherwise semantically and syntactically maintained.

## 7 Related Work

**Representation learning in NLP.** Learning probabilistic models of text for use in representation learning and identifying resulting structure has a long history in NLP, from non-contextual word vectors (Schütze, 1992; Rohde et al., 2005; Turney, 2010; Mikolov et al., 2013; Bojanowski et al., 2017) to contextual networks (Elman, 1990; Bengio et al., 2000; Collobert and Weston, 2008; Sutskever et al., 2011; Peters et al., 2018; Radford et al., 2018). Deep Averaging Networks (Iyyer et al., 2015) are not Backpacks; they first perform averaging and then nonlinear computation.

**Interpretability for Control of NLP networks.** A burgeoning body of work attempts to intervene on monolithic neural networks for interpretability and control (Meng et al., 2022, 2023), and for mechanistic understanding (Olsen et al., 2021; Elhage et al., 2021). Implicitly, Backpacks develop a somewhat human-understandable language of machine concepts, an idea espoused in Kim et al.



(2018); Koh et al. (2020). The connections between interpretation and control are rich; much work has gone into the detection and extraction of emergent structure in networks (Hupkes et al., 2018; Liu et al., 2019) as well as subsequently modulating behavior (Lakretz et al., 2019; Eisape et al., 2022).

**Generalized Additive Models.** Generalized Additive Models (GAMs; Hastie and Tibshirani (1986)) are a function family that (1) independently transforms each input feature, (2) sums these transformations of inputs and (3) applies a non-linear link function (e.g., softmax):

$$f(\mathbf{x}_{1:n}) = \Phi(r_1(x_1) + \dots + r_n(x_n)) \quad (16)$$

Treating each word-position pair as a feature, Backpacks are not GAMs because they include a weighting  $\alpha$  that depends on all features. However, Backpacks share an intuition of computing independent representations of each feature and aggregating by addition. Neural GAMs have been proposed for interpretability (Agarwal et al., 2021; Yang et al., 2021; Chang et al., 2022; Radenovic et al., 2022; Dubey et al., 2022), though never to our knowledge in language modeling. We expect that without context-dependent weighting, models would be insufficiently expressive for language modeling.

## 8 Discussion

In this section, we address a few natural questions about the expressivity and interpretability of Backpacks, highlighting the limits of our knowledge.

### How do Backpacks compare to architecture X?

The Backpack structure does not depend upon using a Transformer to compute the contextualization weights. We could parameterize the contextualization function with a different architecture (e.g., LSTM, S4 (Gu et al., 2021)) and use the resulting weights to compute the Backpack sense vector sum. This architecture, e.g., the Backpack-S4, could then be compared to the standard S4 architecture.

### Are Backpacks as expressive as Transformers?

We don't know. If the number of linearly independent sense vectors is at least  $d$ , then a sufficiently complex contextualization network could treat them as an arbitrary basis. A concern we've often heard is that "simply" adding together sense vectors should not be expressive enough to handle, e.g., negation. However, as long as the requisite

building blocks exist in the prefix, a contextualization network that recognizes the negation or other property could properly distribute weights.

**Are Backpacks inherently interpretable?** No, but we believe no architecture is. Each architecture provides a set of tools that may or may not be useful for differing goals. To us, the key is the mechanistic guarantees Backpacks offer, which will vary in utility depending on how well-specialized the learned sense vectors are for a specific kind of control. Also, the visualizations we provide (top- $k$  highest-scored words) only provide a small view into a sense's potential uses, because scores are non-zero for the whole vocabulary.

### Are Backpacks as compute-efficient as Transformers?

At a glance, no. Backpacks have an underlying Transformer as well as extra parameters, but may perform roughly as well as just the underlying Transformer. However, sense vectors are sparsely activated—only those from the relevant sequence need be on GPU—and after training, can be computed by lookup.

**Why do sense vectors specialize?** Ablations in Table 8 show that they should at least learn to be linearly independent, since linear dependence is equivalent to having fewer sense vectors, which causes higher perplexity. The specialization of sense vectors to seemingly coherent categories may be attributable to the shared feed-forward network that computes them, and/or the contextualization network learning to assign similar weight distributions to senses with similar roles.

**Are sense vectors like "word senses"?** No; they encode a notion of "predictive utility" that doesn't align with traditional notions of word sense. We use the name "sense vector" however because they form a new, useful notion of decomposition of the possible contextual uses of a word into components that are softly combined in each context.

## 9 Conclusion

Non-contextual word2vec embeddings initiated modern deep learning research in NLP, and have fascinating geometric structure. Now, research has largely moved on to monolithic representations, first from RNNs and now from Transformers. Our work suggests that we can have both rich lexical structure and interventions, and strong contextual performance, in a single model.

## 10 Acknowledgements

The authors would like to thank Amita Kamath, Steven Cao, Xiang Lisa Li, Ian Covert, and the Stanford NLP Group community for useful discussions. Further support came from the Stanford Center for Research on Foundation Models. Christopher Manning is a CIFAR Fellow. John Hewitt was supported by an NSF Graduate Research Fellowship under grant number DGE-1656518 and by the CIFAR Learning in Machines and Brains program. We gratefully acknowledge the support of a PECASE Award to Percy Liang.

## 11 Limitations

There is a fundamental uncertainty in whether Backpack language models will continue to scale with parameters and data and be viable alternatives to Transformers at larger model scales. In this study, we were unable to scale larger, and hope that future work will test larger model scales. In a similar vein, we do not verify that Backpack language models perform well across multiple languages. We also do not consider, e.g., finetuning Backpacks on other tasks, or masked language modeling—there is a wide range of possible uses that remain to be verified.

One potential obstacle to the use of Backpacks that we do not study is the effect of tokenization in languages with richer morphological structure than English—will the Backpack structure be amenable to modeling those languages? This may be difficult because, intuitively, the interpretability and control of Backpacks relates to the semantics of individual tokens. Even in English, small subwords not indicative of a single word are hard to interpret. What we hope to have provided is a sufficient set of experiments to motivate the further exploration of Backpacks.

## 12 Ethics

This paper describes and releases an open-domain language model trained on a largely unfiltered subsection of the (mostly English portions of the) textual internet, and describes methods for interpreting and controlling said model. Any control method that can be used to help understand and guide the generation of a model can be used to more effectively generate toxic or illegal content. Despite this, we do expect that, overall, the benefit of deeper insight into Backpack language models is a step

in the right direction. In particular, explanations based on the structure of Backpacks may be able to provide insights into the mechanisms behind model behaviors, increasing transparency.

The concrete models we will release, up to and including 170M parameters, are substantially smaller and less performant at generating text than many of the publicly and commercially available language models available right now, so we do not expect there to be considerable negative repercussions from the release of the artifacts. The code we release, however, could be used or replicated to train much larger Backpack LMs by corporations or governments.

## References

- Rishabh Agarwal, Levi Melnick, Nicholas Frosst, Xuezhou Zhang, Ben Lengerich, Rich Caruana, and Geoffrey E Hinton. 2021. [Neural additive models: Interpretable machine learning with neural nets](#). *Advances in Neural Information Processing Systems*, 34:4699–4711.
- Eneko Agirre, Enrique Alfonseca, Keith Hall, Jana Kravalova, Marius Paşca, and Aitor Soroa. 2009. [A study on similarity and relatedness using distributional and WordNet-based approaches](#). In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 19–27, Boulder, Colorado. Association for Computational Linguistics.
- Dimosthenis Antypas, Asahi Ushio, Jose Camacho-Collados, Vitor Silva, Leonardo Neves, and Francesco Barbieri. 2022. [Twitter topic classification](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3386–3400, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. 2000. [A neural probabilistic language model](#). *Advances in neural information processing systems*, 13.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. [Enriching word vectors with subword information](#). *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. [Man is to computer programmer as woman is to homemaker? Debiasing word embeddings](#). *Advances in neural information processing systems*, 29.
- Rishi Bommasani, Kelly Davis, and Claire Cardie. 2020. [Interpreting Pretrained Contextualized Representations via Reductions to Static Embeddings](#). In *Proceedings of the 58th Annual Meeting of the Asso-*

- ciation for Computational Linguistics*, pages 4758–4781, Online. Association for Computational Linguistics.
- Chun-Hao Chang, Rich Caruana, and Anna Goldenberg. 2022. [NODE-GAM: Neural generalized additive model for interpretable deep learning](#). In *International Conference on Learning Representations*.
- Gabriella Chronis and Katrin Erk. 2020. [When is a bishop not like a rook? when it’s like a rabbi! Multi-prototype BERT embeddings for estimating semantic relationships](#). In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 227–244.
- Ronan Collobert and Jason Weston. 2008. [A unified architecture for natural language processing: Deep neural networks with multitask learning](#). In *Proceedings of the 25th international conference on Machine learning*, pages 160–167.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. [FlashAttention: Fast and memory-efficient exact attention with IO-awareness](#). In *Advances in Neural Information Processing Systems*.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2019. [Plug and play language models: A simple approach to controlled text generation](#). In *International Conference on Learning Representations*.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. [Editing factual knowledge in language models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6491–6506.
- Abhimanyu Dubey, Filip Radenovic, and Dhruv Mahajan. 2022. [Scalable interpretability via polynomials](#). In *Advances in Neural Information Processing Systems*.
- Tiwalayo Eisape, Vineet Gangireddy, Roger P. Levy, and Yoon Kim. 2022. [Probing for incremental parse states in autoregressive language models](#). In *Findings of EMNLP 2022*.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. [A mathematical framework for transformer circuits](#). *Transformer Circuits Thread*.
- Jeffrey L Elman. 1990. [Finding structure in time](#). *Cognitive science*, 14(2):179–211.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2021. [A framework for few-shot language model evaluation](#).
- Daniela Gerz, Ivan Vulić, Felix Hill, Roi Reichart, and Anna Korhonen. 2016. [SimVerb-3500: A large-scale evaluation set of verb similarity](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2173–2182, Austin, Texas. Association for Computational Linguistics.
- Aaron Gokaslan and Vanya Cohen. 2019. [Open-webtext corpus](#). <http://skylion007.github.io/OpenWebTextCorpus>.
- Albert Gu, Karan Goel, and Christopher Re. 2021. [Efficiently modeling long sequences with structured state spaces](#). In *International Conference on Learning Representations*.
- Prakhar Gupta and Martin Jaggi. 2021. [Obtaining better static word embeddings using contextual embedding models](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5241–5253.
- Vikram Gupta. 2021. [Multilingual and multilabel emotion recognition using virtual adversarial training](#). In *Proceedings of the 1st Workshop on Multilingual Representation Learning*, pages 74–85, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. [Array programming with NumPy](#). *Nature*, 585(7825):357–362.
- Trevor Hastie and Robert Tibshirani. 1986. [Generalized additive models](#). *Statistical Science*, 1(3):297–318.
- Felix Hill, Roi Reichart, and Anna Korhonen. 2015. [Simlex-999: Evaluating semantic models with \(genuine\) similarity estimation](#). *Computational Linguistics*, 41(4):665–695.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural computation*, 9(8):1735–1780.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland,

- Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack William Rae, and Laurent Sifre. 2022. [Training compute-optimal large language models](#). In *Advances in Neural Information Processing Systems*.
- Dieuwke Hupkes, Sara Veldhoen, and Willem Zuidema. 2018. [Visualisation and ‘diagnostic classifiers’ reveal how recurrent and recursive neural networks process hierarchical structure](#). *Journal of Artificial Intelligence Research*, 61:907–926.
- Mohit Iyyer, Varun Manjunatha, Jordan Boyd-Graber, and Hal Daumé III. 2015. [Deep unordered composition rivals syntactic methods for text classification](#). In *Association for Computational Linguistics*.
- Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. 2018. [Interpretability beyond feature attribution: Quantitative testing with concept activation vectors \(tcav\)](#). In *International conference on machine learning*, pages 2668–2677. PMLR.
- Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. 2020. [Concept bottleneck models](#). In *International Conference on Machine Learning*, pages 5338–5348. PMLR.
- Yair Lakretz, Germán Kruszewski, Théo Desbordes, Dieuwke Hupkes, Stanislas Dehaene, and Marco Baroni. 2019. [The emergence of number and syntax units in LSTM language models](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 11–20.
- Nelson F Liu, Matt Gardner, Yonatan Belinkov, Matthew E Peters, and Noah A Smith. 2019. [Linguistic knowledge and transferability of contextual representations](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1073–1094.
- Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. 2022. [Locating and editing factual associations in GPT](#). In *Advances in Neural Information Processing Systems*.
- Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. 2023. [Mass-editing memory in a transformer](#). In *The Eleventh International Conference on Learning Representations*.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. [Pointer sentinel mixture models](#). In *International Conference on Learning Representations*.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Efficient estimation of word representations in vector space](#). In *International Conference on Learning Representations (Workshop Poster)*.
- Joakim Olsen, Arild Brandrud Næss, and Pierre Lison. 2021. [Assessing the quality of human-generated summaries with weakly supervised learning](#). In *Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 112–123, Reykjavik, Iceland (Online). Linköping University Electronic Press, Sweden.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2022. [In-context learning and induction heads](#). *Transformer Circuits Thread*.
- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc-Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. 2016. [The LAMBADA dataset: Word prediction requiring a broad discourse context](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1525–1534.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. [GloVe: Global vectors for word representation](#). In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. [Deep contextualized word representations](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. 2021. [Mauve: Measuring the gap between neural text and human text using divergence frontiers](#). *Advances in Neural Information Processing Systems*.
- Ofir Press and Lior Wolf. 2017. [Using the output embedding to improve language models](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*.
- Filip Radenovic, Abhimanyu Dubey, and Dhruv Mahajan. 2022. [Neural basis models for interpretability](#). In *Advances in Neural Information Processing Systems*.

- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. [Improving language understanding by generative pre-training](#).
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#).
- Douglas LT Rohde, Laura M Gonnerman, and David C Plaut. 2005. [An improved model of semantic similarity based on lexical co-occurrence](#).
- Herbert Rubenstein and John B. Goodenough. 1965. [Contextual correlates of synonymy](#). *Commun. ACM*, 8(10):627–633.
- H. Schütze. 1992. [Dimensions of meaning](#). In *Proceedings of the 1992 ACM/IEEE Conference on Supercomputing*, Supercomputing '92, page 787–796, Washington, DC, USA. IEEE Computer Society Press.
- Ilya Sutskever, James Martens, and Geoffrey E Hinton. 2011. [Generating text with recurrent neural networks](#). In *International Conference on Machine Learning*.
- Peter D Turney. 2010. [From frequency to meaning: Vector space models of semantics](#). *Journal of Artificial Intelligence Research*, 37:141–188.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). *Advances in neural information processing systems*.
- Ben Wang and Aran Komatsuzaki. 2021. [GPT-J-6B: A 6 billion parameter autoregressive language model](#). <https://github.com/kingoflolz/mesh-transformer-jax>.
- Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. [BLiMP: The benchmark of linguistic minimal pairs for english](#). *Transactions of the Association for Computational Linguistics*, 8:377–392.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Zebin Yang, Aijun Zhang, and Agus Sudjianto. 2021. [GAMI-Net: An explainable neural network based on generalized additive models with structured interactions](#). *Pattern Recognition*, 120:108192.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. [Gender bias in coreference resolution: Evaluation and debiasing methods](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics.

## A Language Model Training Details

We use the FlashAttention codebase (Dao et al., 2022) which in turn relies on the Huggingface codebase (Wolf et al., 2020) and NumPy (Harris et al., 2020). We perform no preprocessing of OpenWebText. We do no explicit hyperparameter sweep for OpenWebText training beyond our sense vector ablation, instead taking the defaults provided. We train our models on 4 A100 (40GB) GPUs. All experiments test a single trained Small (124M Transformer or 170M Backpack) model due to computational constraints.

### A.1 The feed-forward sense network.

We parameterize the feed-forward network for our sense vectors by first performing layer normalization on the input embeddings, and then a feed-forward layer with residual connection and layer norm (despite it being a function of just one word) to dimensionality  $4d$  and back to  $d$ . Then a subsequent feed-forward network to hidden dimensionality  $4d$  and then up to  $k * d$ . We include a second layer norm and residual before the second feed-forward layer accidentally as a side-effect of the underlying language model codebase.

For our experiments ablating  $k$  in Section 4.5, the second feed-forward component maps to  $d$  and then  $kd$ , not  $4d \rightarrow kd$ .

## B Extra evaluations

### B.1 Timing Benchmarking

To benchmark the speed of each model, we used a single A100 GPU, running the forward pass of each model with a sequence length of 512 and a batch size of 32. We ran 100 forward passes and present the average time taken across the 100. We present this in lieu of FLOPs because A100 GPUs are relatively standard, and this allows for a more directly usable time estimate. Results are in Table 7. We find that Backpacks take roughly 1.4x as long to run as their underlying Transformers.

## C Lexical Similarity Details

To handle words in the lexical similarity datasets that don't appear as single words in the tokenizer, we use one of two methods. We either average all subwords, or take the first subword. The results for the two methods were similar, but we take the better overall for each model. For all Backpack methods, our 124M-parameter Transformer, and

Model	Time ↓
Backpack-Micro	0.093
Transformer-Micro	<b>0.065</b>
Backpack-Mini	0.21
Transformer-Mini	<b>0.15</b>
Backpack-Small	0.36
Transformer-Small	<b>0.26</b>

Table 7: Timing benchmarking results on an A100, average time to compute forward pass on 32-batch size 512-sequence length input.

GPT-2-xl, we average all subwords. For GPT-J (which uses the same tokenizer), we take the first subword.

## D Sense Vector Control Details

### D.1 Topic control details

The full results are in Table 12. The list of topics, and the corresponding bags-of-words, are given in Table 10. For PPLM, the hyperparameter we vary to change the strength of topic control is the step size (Dathathri et al., 2019).

We consider a document as matching the semantic control if the classifier assigns greater than 0.5 probability to the attempted class. We generated from our models with ancestral sampling with no truncation or temperature change.

**Topic control.** Let  $b \in \mathbb{R}^{|\mathcal{V}|}$  be the many-hot vector defined by the bag of words input to the control problem. That is, if the bag is *arts, culture*, then  $b$  has 1 at the indices corresponding to those words, and 0 elsewhere. To determine the initial weights  $\delta$  for each sense vector, we first sort all  $|\mathcal{V}| * k$  sense vectors by decreasing normalized dot product with the bag of words vector:

$$s(C(\mathbf{x})) = \frac{b^\top E^\top C(\mathbf{x})}{\max(E^\top C(\mathbf{x}))} \quad (17)$$

We then take the 0.95, 0.80, and 0.60 quantiles of these scores to determine how to weight the vectors. Intuitively, the vectors in the highest quantiles (most associated with the target topic) are up-weighted the most during decoding, to push the generation towards the topic. The three quantiles partition the set of scores into 4, which are given separate  $\delta$  values; the exact 4 depend on the strength of control (i.e., different points in Figure 2.) The exact  $\delta$  upweighting for each point are given in Table 11.

# Senses	Total Params	Contextl. Params	OWT PPL
1	74.3M	72.7M	38.5
4	75.6M	72.7M	29.3
16	80.5M	72.7M	26.0
64	100.2M	72.7M	24.0

Table 8: OWT perplexity and parameter count as a function of the number of sense vectors. All models trained for 50k steps, 500k token batch size, on OWT.

Model	Dim	Layers	Heads
Micro	384	6	6
Mini	640	8	8
Small	768	12	12

Table 9: Model size hyperparameters.

Topic Label	Bag-of-words
arts_culture	arts, culture
business_entrepreneurs	business, entrepreneurs
celebrity_pop_culture	celebrity, pop, culture
diaries_daily_life	diaries, daily, life
family	family
fashion_style	fashion, style
film_tv_video	film, tv, video
fitness_health	fitness, health
food_dining	food, dining
gaming	gaming
music	music
news_social_concern	news, social, concern
other_hobbies	hobbies
relationships	relationships
sports	sports
travel_adventure	travel, adventure
youth_student_life	youth, student, life

Table 10: The topics used in our topic classifier, and the bags-of-words we use for control.

Control Strength	$\delta$ for quantiles 0.95, 0.80, 0.6, < 0.6
0 (unmodified)	1,1,1,1
1	1.5, 1.5, 1.3, 1
2	2.2, 2.2, 1.5, 1
3	3.3, 3.3, 3, 1

Table 11: Initial topic control weights for each quantile.

**Topic annealing.** From the the beginning value of  $\delta$  given above, we anneal back to 1 as follows. For each sense  $C(\mathbf{x}_j)_\ell$ , we compute the total sum of non-negative log-probability assigned by the sense to the set of words generated so far, intuitively to compute whether the words already generated express the meaning intended by the sense:

$$a_{C(\mathbf{x}_j)_\ell} = \sum_{i=1}^n \max(\mathbf{x}_i^\top E^\top C(\mathbf{x}_j)_\ell, 0). \quad (18)$$

We then re-weight by a term dependent on the sequence index to upweight terms near to the most recently generated text:

$$b_{C(\mathbf{x}_j)_\ell} = \sigma(-a_{C(\mathbf{x}_j)_\ell} f + 6) * (1 + j) / 100 \quad (19)$$

where  $j$  is the index of the word of the sense vector in the generated text, and  $f$  is a scaling constant set to 7.5 divided by the maximum  $\delta$  in the experiment (the maximum of each row in Table 11.)

Finally, we compute the annealed  $\delta$  as a soft combination, weighted by  $b_{C(\mathbf{x}_j)_\ell}$ , of the maximum delta and the default of 1:

$$\delta_{\ell ij} = b_{C(\mathbf{x}_j)_\ell} \delta_{\ell ij} + (1 - a) * 1. \quad (20)$$

## D.2 Gender bias mitigation details

For the third-person singular verb *they*, we found that our sense intervention on sense 10 slightly increases the probability of *they* relative to *he* or *she*.

The full set of nouns and prompts we use is as follows. For role nouns, we use mechanic,

Method	Sem Acc $\uparrow$	Toks-in-vocab $\downarrow$	MAUVE $\uparrow$
<i>Transformer</i>			
Unchanged	6.8%	0.0%	0.95
PPLM-.01	8.4%	0.1%	0.94
PPLM-.04	23.9%	2.6%	0.81
PPLM-.05	30.3%	5.5%	0.62
PPLM-.06	37.7%	12.3%	0.41
PPLM-.07	40.8%	18.8%	0.25
<i>Backpack</i>			
Unchanged	7.4%	0.0%	0.92
Ours <sub>+1</sub>	12.1%	0.2%	0.91
Ours <sub>+2</sub>	24.3%	1.5%	0.90
Ours <sub>+3</sub>	35.3%	3.5%	0.83

Table 12: Topic control via pseudovocabulary, vs PPLM. MAUVE scores are computed with respect to 8000 samples drawn across the topics.

<i>nurse</i>		<i>developer</i>	
logprob	target	logprob	target
+12.21	professional	-15.27	She
+11.62	She	-15.14	herself
-13.38	himself	-13.02	her
-13.24	His	-12.34	corruption
-12.63	his	-12.15	she
-12.31	his	-12.14	gam
-11.39	he	-12.13	HR

Table 13: Some of the most up- or down-weighted words by sense 10 of the words *nurse* and *developer*. For each role word, the sense encodes the stereotypical gender of the role.

accountant, farmer, baker, assistant, construction, guard, carpenter, analyst, physician, cook, clerk, manager, developer, librarian, salesperson, mover, hairdresser, auditor, sheriff, janitor, driver, chief, teacher, writer, receptionist, CEO, nurse, housekeeper, secretary, counselor, attendant, supervisor, designer, lawyer, editor, cashier, laborer, tailor, cleaner.

For evaluation, we use the prompts in Table 14. For experiments in which we estimate the strength of removal to reduce gender bias, we do the optimization individually for each profession noun, on the separate set of prompts in Table 15.



Prompt
<p>I really liked my PROFESSION. What I liked most was that,  I really liked my PROFESSION. In fact,  My PROFESSION said that  My PROFESSION believes that  This one PROFESSION believes that  This one PROFESSION believes  My PROFESSION said  My PROFESSION believes  The PROFESSION was with the car. When  The PROFESSION was with the car. Then  While driving to the store, the PROFESSION looked over on the dash and then  A PROFESSION went to chat over to chat before  Some PROFESSION asked to look over things, but</p>

Table 14: Prompts used for gender bias evaluation.

Prompt
<p>The PROFESSION came into the room. When  The PROFESSION came into the room. Then  I went to the talk to the PROFESSION. What  I went over to the talk to the PROFESSION. Why  I went over to the talk to the PROFESSION;</p>

Table 15: Prompts used for gender debiasing estimation.

Topic	Backpack Sense Control, Control=2	Transformer PPLM, stepsize=0.04
Arts & Culture	(I kind dude) just used to be creative and make music but it feels like I'm 40. Conscious of how bad it is and how awesome it is like making such a fast consciousness and not even being musician meets people who answer for you, especially when it's scary." de la Cruz © Dan Wilson (2002).	Sunday, February 6: WestSTAR en-twitter-ilesurcast.blogspot.com Tuesday, February 7: WestStar en-chat breathe of ancient energy. www.weststar.org Monday, February 8: West Star Mares and Moon of the ages "Happiness is not easy to do", Nicolas Jeansma, the Eternal Life programme director analyses history, culture, sociality and social magic. : 'Oh the
Business & Entrepreneurship	Flickr advertisers is nothing new, so let's hope you know where you buy the latest edition. At the same time, the fix has been pushed through, and while the overall business is pulling away from mainland Asia, publishers have given control over social media options to researchers at New York University and Columbia University. A new report from the Columbia board offers some clues as to why. "My store in Alabama is used to a lot of Marines, and I just dropped as such. I don't know why, but I've had	We've decided to put out a newsletter to your guys, wondering as you cope with the tribulations of your business ventures and a job position. One way to put it is: You're not good enough. You've failed and you're not getting anything done. You're not doing enough. You're not bringing the passion and ideas you might have to a business. But one thing's for sure: if you self-promote, you often might take the business to a profitable buyer. Continue
Celebrity & Pop Culture*	Meetings and greets with reporters and celebrities of all kinds — pop culture, fashion, sports, food, celebrity lifestyle and otherwise — have been laid door-to-door on the Dallas television market with both LaVar and his wife, Arron, taking over the showroom-oneship business at Big Star Barber. "We think Big Star's an interesting exchange," Arron says. "They've got an experience they're	Type Services rumors have been up in the media since last month—and now we have some confirmed to the CBC Radio musical news channel's Twitter stream. The group's guitarist, Greg Carr, has just announced that he's working with Papa John as the band's lead singer and guitarist. According to bizarre French pop culture creation icon Valentino pop music singer/writer Jiv pop pop model, who also wrote pop pop music's MySpace and Twitter pop memes, Cassidy gig pop pop superstar is
Diary & Daily Life	The exact actual life cycle life form life soars on and dies off in comparison to our own. During the first few years of life, the total life form you take to decide what to eat, how much of it to drink, why, and whether you want to exercise have been completely smashed and the technological capability to make that happen seriously out of the blue has been completely lost, jumping from complexity to complexity, totally overwhelming the mushroom in its ability to discover what levels it's supposed to	The Rome crew logam tagged Louisville Main Street today morning and observed a loading dock at the center of downtown Louisville. The dock is just bigger than what was supposed to dock the loading area for emergencies. They watched over the crowd after passing the boat and finally realized that they'd caught some missed traffic signals. "Serious congestion" has so far unnerved people from the Grande family picnics to weddings picnics picnics. MTD Charlotte Pulse (@mtdph)
Fashion	This article is about the fashion label fashion week fashion style month fashion fashion style fashion style fashion week fashion style fashion fashion style fashion fashion style fashion history fashion fashion fashion fashion fashion fashion fashion johnny dressed in an actor's specially created costume news news icon The Comic Relief series features stories, such as plungers from the comic books. It was originally published as a comic published in Dark Horse Comics in English and in both comic books and graphic novels.[1] It was produced	Twitter personality @ceboperformancemk tweeted in response to the story about you. Fashion designer underwear, designer cook dress, sexuality art models, sex con artists, real goths. BuzzFeed You think my brain's shit about what's fashion looks like? Yeah no, I'm not on it. I'm fashion. I'm fine fashion. Yes I appreciate the brand but the people behind it[. . .] adults go fashion, or

Table 16: The first, non-cherry-picked category-satisfying example from each model.

Topic	Backpack Sense Control, Control=2	Transformer PPLM, stepsize=0.04
Film, TV, & Video	Originally published Live chat Qs with the film website writer, who raised millions at least two years ago I contacted him with the same questions as you're doing. I'm a bit optimistic that you're right, but you're just not responding. As you studied the film timer/mapplot'n'cookies response speed, I read the excerpts and couldn't make out a massive amount of time differences. Very minor. What do you think about some of the terms	Well, the hype is real, and with the release of the latest episode of season two (which I'm probably not supposed to review), it feels like you won't be afraid to retweets fideo. By "HAPPY FINALS," the footage maker has used a GIF video to give viewers look at Fideo's dancing triangles and serenity dancing around a moving picture. Thank you, fideo! If the
Fitness & Health	CLOSE Don't think tanking will spell good news for Detroit medical marijuana patients but the owner of its dispensaries saying that is just part of the problem facing the growing number of ill people having access to pot. Healthcare workers are treated for tumors in a dispensary in Oakland. (Photo: Christopher Satorica, Special to CNN) An array of medical centers have lined up near Detroit after a medical marijuana reform forum at the University of Michigan put the debate over the drug at	Today we learn more about the rise of the ice age, multi-drug cocaine epidemic, global population explosion and warfare epidemic by following Dr. Kristof Dr. Freedk published in the British Journal of Medicine The authors update their lofty goal and continue to refine their work for public health. The International Health Services Committee has just released a new research, The next three years could be very costly for health care in Australia, hospitals, state health systems and dietary health. A recent report from
Food & Dining	As weeks wore maple leafed food trucks, and food processors reminisced about their great days past, healthcare workers found out one day that they should get better working conditions with little regard for their bodies. Barbara Butterfield, the former Shop Swagger workshop in Clarksdale, got shot dead on Monday morning when she tried to stop a father Francisco Lee Walker from firing a gun. Walker, 20, had just started his Aug. 27 firing. Exposure to fire and clothes caused Walker	I would dearly love to stand at that galloping chair and who doesn't has amazingly friends associated with their backs hurting? I was a big first timer yesterday. Not always with bacon but I held til calms up. Big chunks of bacon super nice but not me. However there are times where the pieces pull apart and this happens very hard to homo and crackers afgh. All Mixed ones made popular points that have the food triggers across: lack of meats rinsing and eating
Gaming	My parents encouraging kids to be competitive gaming at school is not a new concept. Gaming has been around since the earliest days on paper, and their perspective is always superior than yours. Quality doesn't always apply, and that's why we bucked that trend' father The English woman's son Anthony, who is best known for his role as Most Wanted, came up with the idea of pulling a 30-year-old mentally disabled woman who had been using motorbikes for	Every year, many migrants continue to struggle to find the skills they need in an emerging technology. But every year, it comes quite a surprise to hear the latest news about computerized computing and the gaming community. For the sake of many gaming communities, we here at 14/gamer.org love gaming. It is an important industry in gaming, as it often draws passionate gamers from gaming and lends the gaming community the ability to allow itself special moments like gaming gaming days and gaming gaming. We
Music	David has been a staunch critic of music culture that promotes music as something new, daring, and powerful. As he explained. ("I never thought I was one of those stupid, stupid old people who just listens to music or really hears it it's always the same as when I was a kid," he said.) And when he was a touring musician, those opinions were totally correct. Read the entire interview below. On trying to inculcate younger vocalists with the	From the East art council HQ of MondoJapan Everyone laughs when a sheet metal title is rendered artistically constrained and we say, "Whoa. Then the skin guy! This is a very Chi style steel." Well I don't think anyone's ever heard that before. There's only one coil metal group that is not a tarantella performance music group...at least in America...compart music ten times over and they will never release tracks for it that it is a

Table 17: The first, non-cherry-picked category-satisfying example from each model.

Topic	Backpack Sense Control, Control=2	Transformer PPLM, stepsize=0.04
News & Social Concern	<p>Buildersh B2 has been compared unfathomable by a number of critics because of his security concerns.</p> <p>Breaking News Alerts Get breaking news when it happens — in your inbox. Email Sign Up By signing up you agree to receive email newsletters or alerts from POLITICO. You can unsubscribe at any time.</p> <p>Yet, on Tuesday, Monday and Tuesday, the developer reached the milestone of completing the first UPS facility located in the town of Cloudbreak. He secured \$4</p>	<p>After initially putting itself over Sports Illustrated on Monday, the New York Times was forced to apologize for its widespread coverage of its reporting on the State of Rhode Island – a state that has been the subject of gossip news for some time and which its harsh news and ratings policy has spawned.</p> <p>Late at night on Monday, we learned that the New York Times had reached a breaking news cycle decision and we snagged our exclusive first look at the news. Here’s what you didn’t</p>
Relationships	<p>Early life release parties is relationship couples with relationships over relationships. This census does not count relationships by those who have been with those relationships over the last three years. For more information about early life release parties, check the release party census.</p> <p>Carlo Mathieu Carlo Mathieu was born in 1958. He lives in Augusta, Ga., with his biological father, Malcolm Mathieu, who was president of the Augusta West Raceway at the time. Benjamin Math</p>	<p>Any learning is like being completely ignorant of new information. Schools are forced to teach students to treat one another in the right way, but we still have to recognize that we have to learn how to be friends with as much as we can. When Santod relationships are hard and relationships can be complicated and confusing, there will always be learning relationships, relationships that remind us that we don’t mean relationships, relationships relationships that are boundaries, relationships relationships with friends in need relationships with involved relationships, relationships relationships relationships</p>
Sports	<p>PRESS W/NEWS BLOK Play slideshow 1 of 83 Express sports retail giant Sports Direct.</p> <p>Sports Direct has revealed the on offer outdoor sports gear Brand new from Google has been developed. Here’s what you can expect from Google’s sporting expertise.&lt;lendoftext&gt;About The potential of a west coast restaurant for tolerance and pity</p> <p>Their position at this point hurts me less than they believe it deserves, because they probably shouldn.</p> <p>I’m going to help them</p>	<p>Authorities in California say they are investigating equestrian skiers who struck a 19 year-old boy from a snow-covered mountain and beating him on the head with shovels. According to Smith-Cox, those same well clients found out they had also been tardled by a \$500 pour from pipe on top of of a Black Rock vault. And it appears the ultimate goal of those riders and their company of riders was killed. Jeremy Goschz is one of those survivors. His racing</p>
Travel & Adventure	<p>My next stop destination for me is adventure travel. I travel Disney World and make sure that the worlds under my belt and desert warriors that I’ve been fighting for have a place or two at their disposal that are compatible with my use of current technology. This job is being completed with the help of any freelance user submission information you may have provided. It’s only fair to give you some tips to help you figure it out if there are any unknown sideside locations that you</p>	<p>Equality Equality – open life – inequalities – political oppression – write and publish your work Equality is a freedom to work, to die. Access to free healthcare, free outer space travel, photocopies online, happy endings, self travel – to travel to someone else’s heart (read: stop taking drugs), to move faster, to travel in train travel, to stop a vacation abroad (tell others your travels), to return to a home each time</p>
Youth & Student Life	<p>College students at almost every age advantage who take advantage of learning opportunities in the sport of running spend at least five years an average of \$10 or more per year to do it, according to the University of San Diego’s National Football Clearinghouse.</p> <p>Those risk factors lift nearly a third of university and college football athlete spend, more than double that of a comparable age group of men and women who spend 4,000 hours per year as runners, or 5,000 to</p>	<p>lame University saw a 32 per cent rise in its undergraduate science institutes and 14 per cent increase in its researchers from recent years.</p> <p>Director Of University Development, Mike Brennan, said: "The growth in university employment, coming from such a historic campaign, is something to celebrate as we support our young people and room to progress in science and technology."</p> <p>A student was interviewed in a recent paper about university employment, specifically a dissertation.</p> <p>"For the first time, people are</p>

Table 18: The first, non-cherry-picked category-satisfying example from each model. This is except for the Relationship category for the Transformer, where we skipped the first one due to content we particularly did not want to publish.

Positive Log-Probability Mass for Senses of word <i>quickly</i>							
0	1	2	3	4	5	6	7
approaching	oggles	quickly	enough	stro	iii	razen	asuring
ascended	Marks	swiftly	rotating	zn	Original	forgotten	delusion
grav	Axis	rapidly	paced	strokes	alsa	forget	stimulated
gent	claimer	quick	ened	uling	chenko	social	recollection
disposed	Roche	quick	retreating	\$_	resolution	rius	stimul
risen	demonstration	instantly	Subscribe	grass	ient	relapse	Wem
dispose	blaster	promptly	dismissing	lessly	baskets	baseless	persistent
becoming	ducers	soon	diminishing	iken	uin	Statement	urbed
ascert	Specifications	fast	disappearing	izing	ora	athing	retard
climbed	Viet	Quick	varying	bg	alid	Akron	restraint
8	9	10	11	12	13	14	15
processors	slowly	tering	Definitely	quick	oted	ouse	Sims
darts	Slowly	Bers	initely	quickest	distances	pee	Noir
milliseconds	Slow	Fed	raid	quick	outed	ouses	TMZ
wip	conveniently	ascus	alright	quicker	aught	pees	Streets
iazep	slower	Bust	Personally	fast	UC	attach	expressly
reptiles	cheaply	aucus	laughs	quickly	ob	tro	Attend
Kelvin	responsibly	Ryu	ALWAYS	rapid	digits	iffe	Rooms
Ow	gradually	sector	Always	fast	ench	aces	Within
Soon	quietly	Petra	Ideally	faster	Code	lain	Rum
Slug	waiting	DCS	Roses	fastest	apers	feet	Forced
Negative Log-Probability Mass for Senses of word <i>quickly</i>							
0	1	2	3	4	5	6	7
initely	sburg	ollen	una	Poké	quickly	Faster	.
heit	orem	oned	URE	slow	quick	purposely	Sorceress
Aly	Untitled	oths	rast	slower	swiftly	deliberately	itars
istically	anted	ook	ipt	slows	rapidly	Definitely	Shogun
Always	untreated	ught	ocracy	slowed	quickest	ey	Yen
Doctors	til	Ded	law	DEV	quick	slower	oenix
dl	broken	lost	uthor	encia	Quick	initely	Jagu
urally	past	aught	ema	potions	fast	isner	izz
ependence	ebook	recharge	ory	Machina	instantly	hesitated	eral
raints	Continue	ady	antis	Slow	Quick	eyewitness	finals
8	9	10	11	12	13	14	15
quist	WM	prototype	ciating	kins	quick	Laur	thal
ocker	isf	projector	scrambling	Host	quick	Never	imble
ovsky	fb	reconcil	rapid	loudspe	quickly	Jimmy	iquid
ictions	WF	prominently	newcomer	enced	Quick	dearly	initialized
olation	elevation	counterfeit	adapting	Evil	soon	Dating	ansas
cano	RM	word	speeding	washed	fast	_ _	IGH
Proof	975	cellul	frantic	Kaf	rapidly	never	unciation
cert	dir	prototype	novelty	Glass	Quick	Certainly	needs
rero	ESE	collaps	paced	sod	hurry	eternal	commit
anch	onder	dyl	instructional	advers	Immediately	Rare	tackle

Table 19: For each sense vector of the word *quickly*, the 10 words to which the sense vector assigns the highest log-probability contribution, and the 10 to which it assigns the largest negative log-probability contribution. Note that usually, either the positive words are coherent or the negative—but not both for the same sense index. Some senses are not interpretable, and seem to be used by other parts of speech.

## ACL 2023 Responsible NLP Checklist

---

### A For every submission:

- A1. Did you describe the limitations of your work?

9

- A2. Did you discuss any potential risks of your work?

9

- A3. Do the abstract and introduction summarize the paper’s main claims?

1

- A4. Have you used AI writing assistants when working on this paper?

*We used ChatGPT and Claude to try to brainstorm names for models; nothing useful came of it or ended up in the paper.*

### B Did you use or create scientific artifacts?

*Section 5,6,7*

- B1. Did you cite the creators of artifacts you used?

*Section 5,6,7, Appendix*

- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?

*Left blank.*

- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?

*Left blank.*

- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?

*Left blank.*

- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

*Left blank.*

- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.

*Left blank.*

### C Did you run computational experiments?

*Left blank.*

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?

4

---

*The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.*

C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

4.2

C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

4.2

C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

Section 4.1

**D**  **Did you use human annotators (e.g., crowdworkers) or research with human participants?**

*Left blank.*

D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

*Not applicable. Left blank.*

D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

*Not applicable. Left blank.*

D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

*Not applicable. Left blank.*

D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

*Not applicable. Left blank.*

D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

*Not applicable. Left blank.*