

RellE	FeatID	Interpreted Function
<b>Comparison: 128M - 1B</b>		
<i>128M specific</i>		
0.19	5667	-
0.28	14250	Detects token <i>-ese</i> at the end of a word
0.29	440	Detects token <i>-ara</i> at the end of a name, promotes possession or verbs
<i>128M-1B shared</i>		
0.38	8636	-
0.43	3164	Detects <i>-us</i> ending, often for a Latin origin single noun and promotes verb <i>is</i>
0.52	12683	Detects the noun <i>analysis</i>
0.53	1749	Detects irregular plural noun <i>people</i>
0.56	5032	Detects irregular plural noun <i>men</i> , promotes EOS, conjunction, or verbs
0.57	7072	-
<i>1B specific</i>		
0.71	15882	Detects singular <i>man</i> , promotes preposition completion
0.83	4118	Detects singular <i>woman</i> , not necessarily as a subject
0.89	6381	-
0.91	10069	Detects nouns that end with <i>-ists</i> and promotes plural verb completion or prepositions
0.92	14897	Detects nouns that end with <i>-ans</i>
0.94	16118	Detects words ending with <i>-ias</i>
0.95	8757	-
1.00	3811	Detects regular plural nouns and promotes conjunction or prepositions
1.00	7483	Detects singular nouns preceded by <i>this</i>
<b>Comparison: 1B - 4B</b>		
<i>1B specific</i>		
0.00	12677	Detects regular plural nouns that refer to groups of people and promotes plural verb completion or prepositions
0.10	5778	Detects <i>man</i> starting words but promotes multi-token word completions (, <i>-hood</i> , <i>-ned</i> , <i>-hattan</i> )
0.17	1440	Detects words containing the mid-token <i>-es</i> and promotes medical term completions (, <i>-esophagus</i> )
0.28	3737	Detects singular <i>woman</i> , not necessarily as a subject
<i>1B-4B shared</i>		
0.37	12685	Detects nouns that end with <i>-ans</i>
0.48	7616	Detects nouns that end with <i>-ists</i> and promotes plural verb completion or prepositions
0.53	14814	Detects singular nouns preceded by <i>this</i> in front of them, and promotes singular verb conjugation
0.68	11799	Detects regular plural nouns that refer to objects/science concepts and promotes plural verb completion or prepositions
<i>4B specific</i>		
0.82	9385	-
0.85	744	-
0.88	14210	Detects regular plural nouns
1.00	13102	-
1.00	1868	Detects punctuations and newline to promote BOS words
1.00	4050	Detects nouns that are preceded by plural quantifiers (, <i>most</i> , <i>many</i> , <i>majority of</i> , <i>some</i> )
1.00	9326	Detects plural regular nouns, promotes plural conjugated verbs
1.00	9414	Comma detector
1.00	11088	Detects the final token of first names
1.00	14546	Detects HTML/code-related regular plural object nouns
<b>Comparison: 4B - 286B</b>		
<i>4B specific</i>		
0.00	4368	-
0.00	2307	Detects regular plural nouns that refer to science concepts and promotes plural verb completion or prepositions
0.14	479	-
0.19	680	-
0.27	13452	Detects (regular and irregular) plural nouns, promotes plural verb completion
0.27	10514	Detects regular plural nouns that can be followed up with <i>themselves</i>
<i>4B-286B shared</i>		
0.37	15084	Detects (regular and irregular) plural nouns that refer to groups of people, promotes plural verb completion
0.56	5268	-
0.58	5129	-
<i>286B specific</i>		
0.79	14815	-
0.85	6511	-
0.89	5588	Detects newlines
0.92	12003	-
0.98	13244	Detects first names that are not followed up a last name
0.98	12108	Detects HTML/code-related plural object nouns
1.00	2139	Detects a larger variety of prepositions and complementizers (, <i>by</i> , <i>from</i> , <i>due</i> , <i>with</i> , <i>concerning</i> )
1.00	5138	Detects last token of multi-token first names, promotes last names

**Table 1: 2-way L1-Sparsity Crosscoder Annotation for Pythia-1B.** Each block is one pairwise comparison. is sorted from 0.00 to 1.00, where  $\downarrow 0.3$  gets attributed to the first checkpoint;  $\downarrow 0.7$  to second; shared otherwise). Interpreted function gives a description if a linguistic role was detected, “—” otherwise. Pairwise comparisons reveal finer-grained feature shifts from one checkpoint to another, but cannot assess persistence like triplet analyses. This shows that early checkpoints (, 128M) capture low-level lexical and morphological patterns, while slightly further trained ones (, 1B) detect slightly more abstract patterns, such as irregular plurals.