

RelIE	FeatID	Interpreted Function
1B-4B shared		
[0.53, 0.33, 0.15]	1067	Detects subtoken <i>-ans</i> typically in names
[0.41, 0.39, 0.20]	941	Detects plural nouns that are art-related professions
[0.45, 0.41, 0.14]	4897	Detects plural nouns that end with <i>-ists</i> , (, <i>protagonist, capitalist, pharmacist</i>)
[0.32, 0.43, 0.25]	15204	Detects singular nouns found in technical discourse (, <i>method, function, guide, recipe</i>) preceeded by the word "This"
1B-286B shared		
[0.55, 0.10, 0.34]	7489	Detects singular <i>woman</i> noun
[0.52, 0.03, 0.45]	1641	Detects newlines
[0.52, 0.01, 0.46]	3852	Detects singular <i>man</i> noun
4B specific		
[0.00, 1.00, 0.00]	15556	Detects a full stop and promotes connection words or newlines
[0.00, 1.00, 0.00]	11274	Multi-word noun or compound noun detector
[0.00, 0.99, 0.01]	8318	Detects regular plural nouns
[0.00, 0.96, 0.04]	10020	-
[0.11, 0.69, 0.20]	15950	Detects regular plural nouns
[0.02, 0.68, 0.30]	10523	Detects plural nouns mostly depicting humans (, <i>people, students, bloggers</i>)
[0.11, 0.62, 0.26]	15118	-
4B-286B shared		
[0.01, 0.30, 0.69]	11987	-
286B specific		
[0.00, 0.00, 1.00]	15323	Detects plural nouns found in technical discourse
[0.08, 0.15, 0.77]	14228	Multi word named entity detector (proper nouns, locations etc.)
[0.00, 0.00, 1.00]	15027	Activates on last token of capitalized names (person, location etc.)
[0.00, 0.18, 0.82]	6746	Detects deverbal nouns / nominalizations, abstract/eventive nouns formed from verbs
[0.00, 0.10, 0.90]	5317	-
[0.01, 0.23, 0.76]	14629	Newline detector
[0.10, 0.08, 0.82]	13117	Newline detector
[0.00, 0.00, 1.00]	15129	First name detector
[0.00, 0.01, 0.99]	14623	Detects prepositions

Table 1: 3-way L1-Sparsity Crosscoder Annotation for Pythia-1B — Comparison 1B 4B 286B. shows 3-way one-versus-all vector; Interpreted Function provides a description if a linguistic role was detected, and “—” otherwise. Rows are grouped by checkpoint specificity according to : features dominated by one checkpoint (1B, 4B, 286B specific); pairwise shared features (1B-4B, 1B-286B, 4B-286B shared); and shared across all (1B-4B-286B shared). A missing group means no such features found in the top-10 IE features of all checkpoints. -based triplet comparisons reveal that earlier checkpoints (, 1B and 4B) primarily detect low-level lexical and morphological patterns such as suffixes and irregular plurals, whereas later checkpoints (, 286B) increasingly specialize in higher-level syntactic and semantic functions, including named entity, nominalization, and technical discourse related noun detection.