| RelIE | FeatID | Interpreted Function | Languages |
|---|---|---|---|
| **Comparison: 550M  6B** | | | |
| *550M specific* | | | |
| 0.00 | 8760 | Detects administrative/government-related nouns | fra |
| 0.00 | 14133 | Detects nouns and verbs that convey key actions, entities, or ideas in a sentence | eng |
| 0.05 | 12275 | Detects conjunction token *et* | fra |
| 0.10 | 8341 | Detects subtoken *et* typically conjunction but also for *et al.* | fra |
| 0.20 | 15852 | Detects head nouns and their modifiers that signal prominent participants or components | fra |
| 0.22 | 14697 | Detects plural nouns, promotes plural verb conjugation and *who* pronoun | eng,fra,spa |
| *550M-6B shared* | | | |
| 0.36 | 1474 | Detects plural French articles (, *les, nos, certains*), promotes plural single token noun completions | fra |
| 0.39 | 8223 | Promotes *-age* completion for nouns | eng,fra,spa |
| 0.44 | 7882 | Detects English relative pronoun *that*, promotes pronoun follow ups | eng |
| 0.69 | 14645 | Detects *Ev* at BOS, to be completed with French adverbs or nouns | fra |
| *6B specific* | | | |
| 0.95 | 12523 | - | - |
| 0.96 | 10853 | Detects noun and nominal expressions representing abstract entities, events, or processes | eng,fra,por,spa |
| 0.97 | 2189 | Detects multi-word nouns (, compound nouns or with adjectives) | eng,fra,por,spa |
| 1.00 | 9386 | - | - |
| 1.00 | 9813 | *who* and *that* detector | arb,eng,fra,por,spa |
| 1.00 | 10337 | Punctuation and newline detector | arb,eng,fra,por,spa |
| 1.00 | 14067 | - | - |
| 1.00 | 15428 | Detects verbs with the concept to like/love/appreciate | eng,fra,hin*,por,spa |
| **Comparison: 6B  55B** | | | |
| *6B specific* | | | |
| 0.00 | 11469 | Detects punctuation and newline | eng,fra,spa |
| 0.17 | 942 | Verbs that depict dynamic, agentive actions | eng |
| 0.20 | 10311 | Detects *ev* or *Ev* tokens to be completed with french adverbs or nouns | fra |
| 0.28 | 15632 | Detects verbs with the concept to like/love/appreciate | arb*,eng,fra,por,spa |
| *6B-55B shared* | | | |
| 0.31 | 11920 | Detects head of noun phrases denoting concrete or informational entities (,*data, system, text*) | eng,fra,por,spa |
| 0.32 | 12000 | Detects noun that depict occupational or social roles like researchers, engineers, physician, and journalists | arb*,eng,fra,hin*,por,spa |
| 0.32 | 2792 | - | - |
| 0.38 | 14748 | - | - |
| 0.40 | 5763 | Detects the concept boss in different languages (, *chef, boss, jefe*) probably because it is a common noun in the IE dataset | eng,fra,spa |
| 0.56 | 9817 | Detects verbs and promotes preposition/conjunction/punctuation | eng,fra,por,spa |
| 0.56 | 425 | Detects relative pronoun and prepositions (, *that, que, at, of, de*) | eng,fra,spa |
| 0.59 | 4863 | Detects relative pronoun *that* and promotes verb/pronoun completions | eng |
| *55B specific* | | | |
| 0.79 | 5345 | Nouns and verbs related to consultants and consulting | arb*,eng,fra,por,spa |
| **Comparison: 55B  341B** | | | |
| *55B specific* | | | |
| 0.27 | 15249 | Nouns and verbs related to consultants and consulting | arb*,eng,fra,por,spa |
| 0.27 | 12734 | - | - |
| *55B-341B shared* | | | |
| 0.35 | 11458 | Detects *that* and promotes verb or pronoun completion in English | eng |
| 0.39 | 11920 | Detects nouns to be completed by *de, of* | eng,fra,por,spa |
| 0.39 | 7339 | Determiners and quantifiers in noun phrases, such as articles (, *a, os*), possessives ( *our*), and universal quantifiers ( *every, todo, cada*) | eng,por,spa |
| 0.44 | 6063 | Detects the concept boss in different languages (, *chef, boss, jefe*) probably because it is a common noun in the IE dataset and promotes *of* (, *de, do, du*) completions | eng,fra,por,spa |
| 0.48 | 15083 | Relative pronouns and the syntactic material inside relative clauses | eng,fra,por,spa |
| 0.54 | 10325 | Detects *ev* or *Ev* tokens to be completed with french adverbs or nouns | fra,por,spa |
| 0.55 | 425 | Detects relative pronouns/subordinators (, *that, que, qui, which, who, où*) to introduce a new clause; also activates on the verbs inside the subordinate clause | arb*,eng,fra,por,spa |
| 0.57 | 8729 | - | - |
| *341B specific* | | | |
| 1.00 | 794 | - | - |
| 1.00 | 7419 | Detects newlines in different languages | arb,fra,hin,por,spa |
| 1.00 | 7806 | - | - |
| 1.00 | 13276 | - | - |
| 1.00 | 13404 | Sentence-boundary detector through punctuation and other delimiters | arb,fra,por,spa,zh |

Table 1: **2-way L1-Sparsity Crosscoder CLAMS French/English Annotation for BLOOM-1B.** Each block is one pairwise comparison.  is sorted from 0.00 to 1.00, where ¡ 0.3 gets attributed to the first checkpoint; ¿ 0.7 to second; shared otherwise). Interpreted function gives a description if a linguistic role was detected, "–" otherwise. Languages lists which languages the feature highly activates on, * means that the activation was relatively less common. While earlier checkpoints (, 550M) capture language specific low-level function words, later checkpoints (, 55B and 341B) increasingly share such features across languages.