

Layer	Lang	Feature ID	Top Tokens
layers.0.mlp	Korean	87423	kim, velt, Office, PMC, orney, ordinances, agre, %/, agony, office
layers.4.mlp	Korean	99278	Kim, Seoul, Jae, Kim, .kr, Samsung, Korean, Korea, Shin, Koreans
layers.6.mlp	Korean	109962	Highlands, ussels, iteur, tails, Ire, ativ, \$sql, plant, jsp, beginnings
layers.7.mlp	Korean	115377	vain, ane, ANE, lett, att, anes, roc, 区, igne, olina
layers.8.mlp	Korean	94922	ante, iente, innie, ensed, λλη, ssh, flips, uş, Liter, itar
layers.9.mlp	Korean	31611	beck, acles, ataka, batt, -wrap, etat, pty, umn, Kear, chat
layers.9.mlp	Korean	53931	iosa, gli, chet, conì, Giov, 89, adero, ewolf, worm, edit
layers.10.mlp	Korean	96918	-Feb, partial, juan, pur, ipro, capitals, URATION, 膜, [, experiences
layers.10.mlp	Korean	125019	当, Mori, Petty, Yug, Touch, _che, bare, VIC, .jp, omi
layers.11.mlp	Korean	64252	viz, avour, aura, 今日, 琴, UPDATED, apl, intend, intention, ITT
layers.11.mlp	Korean	102511	r, 으로, _FR, aar, l, war, 은, Protocol, айд, 이라
layers.12.mlp	Korean	27775	ptime, epochs, ascending, located, plays, inf, idual, Microsystems, 舖, layers
layers.12.mlp	Korean	67845	Kim, Korea, Korean, Seoul, Park, Kim, Hyundai, kim, Koreans, Je
layers.13.mlp	Korean	78512	plays, 바람, play, Maar, 편, Boost, maz, robots, antics, 帖
layers.14.mlp	Korean	87496	상, 령, 선, 이상, 나, 신, 증, 등, 한, 왕
layers.14.mlp	Korean	107903	asc, ey, anz, erc, acic, eh, ans, eyi, ekl, erv
layers.15.mlp	Korean	17864	이야기, 그것, 있었다, 없었다, 사람들이, 인정, 행동, 사람은, 이야, 그를
layers.15.mlp	Korean	41829	척, 찰, 리로, 리카, 림, ㄹ, 어나, 린이, 니아, 택
layers.15.mlp	Korean	57880	?, ? , ?", ?", ?" , ?', ?(, ?, , ?', ?:
layers.15.mlp	Korean	60936	통, 작, 소, 지, 그, 기, 판, 누, 빙, 식
layers.15.mlp	Korean	114116	LETE, Passive, 설명, flashlight, 게임, careful, REVIEW, ▼, 中文, 问
layers.15.mlp	Korean	131027	,, reckon, örper, 림, ., 간, 모습, ,',, 손, лива