

1. Top-k, medium-k and bottom-k phrases in single/multi-words rank list for each dataset, (k = 30.)

Dataset	DBLP.300K			
	Single-word		Multi-words	
top-k	collins hebrew basque paul sony persian berlin siemens osgi rome gibbs kim september michigan thomas	kripke macintosh tiger nokia cray coral ramsey swiss gödel cayley sdram eiffel michael india riemannian	simulated annealing quality assurance neural network natural language augmented reality wavelet transform mixed reality wireless lan markov chains elliptic curve turing machines relevance feedback vector quantization timed automata decision tree	edit distance anomaly detection gene expression belief propagation feature selection matrix multiplication markov chain finite state machines virtual reality neural networks hash functions maximum entropy reinforcement learning concurrency control genetic algorithm
medium-k	webserver normalised calibration jb clfsr seats conceptually circles friendliness poly equivalents extractors basketball multitask tailor	centers environments sign leadership episode reorganization equivalences mismatch vector compatibility juntas trustees decommitment conical facile	up to 30 be reviewed segmentation via better than more general n point still lack based on adaptive researchers who to verify some interesting a preprocessor propose a heuristic paper describes how lossy compression of	an aggregation approach to the problem method to identify develop an efficient existing ones techniques for efficient demonstrate the performance can be accessed reason is that turn out to be major drawback of approach to reduce algorithm enables has the advantage reaction time
bottom-k	22 could taken used 91 viz 05 42 come whole 67 except 72 800 120	34 beforehand 000 160 placed seems 84 02 150 seem cannot 81 concerning indicates indicate	system which has with each other in this paper we describe the are interested in the to the identification of be used in the can be used by find all the to the characteristics of have designed a is that it can as described in as an example to show that they is well known that	new approach for the be one of the with a series this paper we show of our approach on in several different be adapted to the is based on the use has been done to an important role in the been applied to the is a need for that this approach is we discuss some of been used in the

Dataset	YELP.100K			
	Single-word		Multi-words	
top-k	alas tiki cuban coach verde latin village capitol antipasto queso wildflower pierre dental gilbert northern	bourbon monterey cherry mcdonalds manhattan john railroad sprinkles hurricane zinburger national macy's brie hula safeway	sea salt papa johns barrio cafe la fitness cotton candy humble pie smoked salmon crab puffs chopped liver peanut butter bell pepper yellow tail pumpkin porter chile relleno pork chop	soda fountain daily dose fettuccine alfredo urgent care del rey vanilla bean cheesecake factory fountain hills collard greens bikram yoga golden corral lettuce wraps dairy queen los betos wild boar
medium-k	extensive sweet.the envious babysitting somebody's elves good.for restless highs clean.i sparkle curves regrettably swore fryer	intrigued chair providers promote existence round choc relive joints share uninspired memorabilia calendars onions actual	to feed was incredibly we were served an important and the selection want a place sat near everything is cooked every inch of i've spent easy on the section of the runs from keeping them you actually get	who was visiting let him i did enjoy bothered to i haven't was so cool i really feel and the noodles was also nice had the grilled were also pretty rice & to the chandler removed from the within 5
bottom-k	105 62 87 1200 1960 97 180 \ 07 83 11.00 98 05 seeming	concerning 69 2002 1999 indicates 59 72 arent 04 causes 06 01 1950 shed 350	you live in the and the place has to speak to a and the salsa was be a bit more i guess if you is located in the on the list of i was glad i is a great place for a is very close to and the cheese was a visit to the this location has a was the best part of	had a pretty good in the center of in a place like i really enjoy this i was glad to my husband is a for lunch on a out to be a have only been to have been coming to on one side of is a place you have a variety of were seated in the had to wait about

2. A table includes number of qualified phrases in each dataset, average number of phrases in each sentence.

	DBLP.300K	YELP.100K
Number of qualified phrases	956,355	1,560,970
Average number of phrases per sentence	0.677733	0.866286

3. Print several clusters and 20 words in each cluster. Describe Clustering method you are using and number centers you've set.

I use `nlk.cluster.kmeans.KMeansClusterer` and set the number of centers = 20.

DBLP.300K	
Cluster 0 (description: Computer Science) _object-oriented_design_ object-_relational_database_systems_ _backward-compatible_ metaprogramming _database_engine_ _data-flow_ _testing_methodology_ _finite_state_machine_ _net_framework_ _database_schema_ _design_tool_ _abstract_state_machines_ _abstraction_levels_ _source-code_ bigraphical _analysis_tools_ _fault_injection_tool_ _verification_environment_ _machine_language_ _functional_programming_languages_ 	Cluster 1 (description: Natural Language Processing) _noun-phrase_ _word_sense_ _opinion_mining_ _query-focused _cross_lingual_ _clustering_approaches_ _translation_examples_ _semantic_context_ _delimiter_ _text_corpora_ _language_translation_ _hmm_based_ _linguistic_analysis_ vocabularies _parallel_corpus_ _computer-assisted_translation_ _textual_data_ enriching _automatically_extracts_ transliterations
Cluster 2 (description: Machine Learning) _optimization_algorithms_ _multiagent_reinforcement_learning_ _hybrid_methods_ _sequence_alignments_ annealing em-like _based_adaptive_ _semi-markov_ _continuous_state_ _approximate_inference_ _multiagent_learning_ _multiclass_ _nearest_neighbor_queries_ _hierarchical_bayesian_ _propagation_method_ filter _curse_of_dimensionality_ _clustering_techniques_ _feature_weights_ mlps 	Cluster 3 (description: Computer Architecture) _fpga_devices_ _processor_arrays_ dma 8-node _heterogeneous_multiprocessors_ _back_end_ microprogrammed synthesizer _soc_designs_ _based_systems_ _workload_characterization_ _bus-based_ micro-processor _dynamically-scheduled_ _parallel_machine_ _super-scalar_ _highly-parallel_ non-pipelined _single_processor_ _tasks_scheduling_

YELP.100K	
Cluster 0 (description: Time) _friday_afternoon_ _pretty_late_ morning's _late_night_dinner_ _sunday_nights_ _black_friday_ _tuesday_morning_ _friends_birthday_ _winter_ _soft_opening_ _mid-afternoon_ after-hours _mothers_day_ _business_lunch_ _busy_nights_ _weekday_night_ weekends _open_til_ _weekend_night_ workday	Cluster 1 (description: Feelings) surprises _forgettable_ _low_rating_ unpredictable _letdown_ horrible _dislike_ critical _lame_ _decent_breakfast_ underrated uninspired wasteful _five-star_ _memorable_experience_ _highly_doubt_ _food_poisoning_ _high_prices_ _horrible_food_ _biggest_issue_
Cluster 5 (description: Desserts and Drinks) _sour_cherry_ pastry _chocolate_torte_ _red_bean_paste_ _homemade_ice_cream_ smooth milkshake caramels _kettle_corn_ espresso tea _chocolate_milk_ _shortcake_ _pudding_ _strawberry_sauce_ _raspberry_sorbet_ _lemon_tart_ _gingerbread_ _godiva_ _waffle_cone_	Cluster 10 (description: Opinions) leftovers _casual_dining_ _average_prices_ _fairly_decent_ _decent_quality_ _thumbs_up_ _yummy_food_ _awesome_ _decently_priced_ _friendly_atmosphere_ _dirt_cheap_ _attentive_service_ _reasonable_prices_ _cheap_meal_ _outstanding_food_ _neighborhood_place_ _casual_atmosphere_ _decent_service_ _tastes_amazing_ well-versed

4. A parameter study on AutoPhrase and Clustering, for example, by changing HIGHLIGHT_MULTI, HIGHLIGHT_SINGLE from 0 to 1 in 0.2 increments in phrasal_segmentation.sh you can get different number of phrases in corpus.

(1) Draw the number of phrases versus HIGHLIGHT_THRESHOLD curve for both multi-words and single word phrases returned by AutoPhrase.

Fix HIGHLIGHT_MULTI to 1.0 and adjust HIGHLIGHT_SINGLE from 0.0 to 1.0 in 0.2 increments.

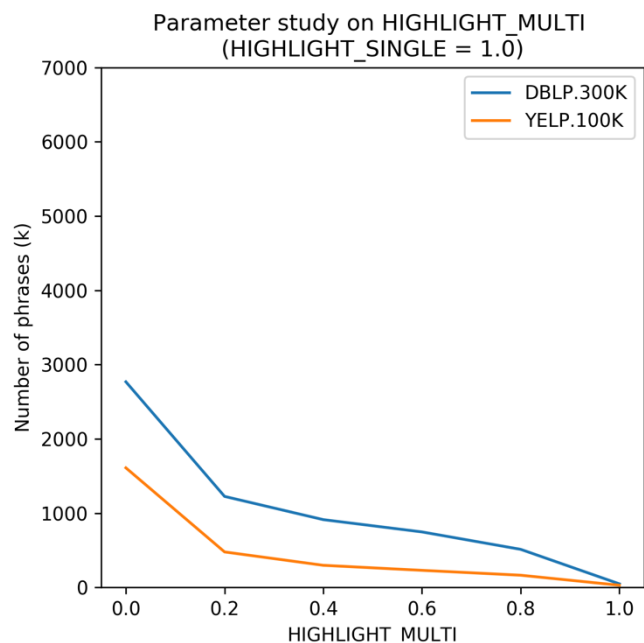
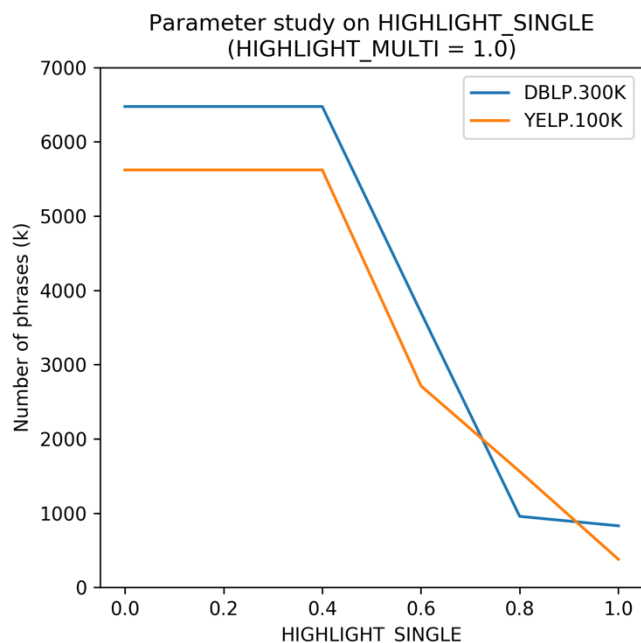
DBLP.300K						
HIGHLIGHT_SINGLE	0.0	0.2	0.4	0.6	0.8	1.0
HIGHLIGHT_MULTI	1.0	1.0	1.0	1.0	1.0	1.0
Number of Phrases	6,474,506	6,474,506	6,474,506	3,700,344	956,355	829,164

YELP.100K						
HIGHLIGHT_SINGLE	0.0	0.2	0.4	0.6	0.8	1.0
HIGHLIGHT_MULTI	1.0	1.0	1.0	1.0	1.0	1.0
Number of Phrases	5,620,537	5,620,537	5,620,537	2,712,355	1,560,970	378,531

Fix HIGHLIGHT_SINGLE to 1.0 and adjust HIGHLIGHT_MULTI from 0.0 to 1.0 in 0.2 increments.

DBLP.300K						
HIGHLIGHT_SINGLE	1.0	1.0	1.0	1.0	1.0	1.0
HIGHLIGHT_MULTI	0.0	0.2	0.4	0.6	0.8	1.0
Number of Phrases	2,767,982	1,224,934	912,284	746,795	512,180	47,565

YELP.100K						
HIGHLIGHT_SINGLE	1.0	1.0	1.0	1.0	1.0	1.0
HIGHLIGHT_MULTI	0.0	0.2	0.4	0.6	0.8	1.0
Number of Phrases	1,609,528	477,466	297,388	229,658	163,845	28,895



(2) By setting number of centers, you can get phrase clusters in different granularity. Show some representative clusters and 10 words in each cluster for different granularity. (e.g. k = 5, 10, 25)

DBLP.300K		
k= 5	k = 10	k = 25
Cluster 0 (description:) theweb domain-specific _relational_data_base_ _knowledge_structures_ _egee_ _distributed_hash_tables_ _language-based_ _separation_of_duty_ neues konturen	Cluster 1 (description:) _web_communities_ _inverted_file_ _precision/recall_ pedigree _holistic_twig_ _text_classifier_ _learning_phase_ _keyphrase_extraction_ _quad_tree_ _bilingual_dictionaries_	Cluster 0 (description: Computer Science) _view_integration_ _escher_ _state_space_analysis_ _kit_ marte _standard_uml_ erts _owl-s_ _domain_specific_language_ sbql
Cluster 1 (description:) _hemoglobin_ garch _single_frame_ bias tagging btfs _document_space_ inflection _hyper_-spectral _newton_method_	Cluster 3 (description:) _reed-solomon_ _instruction-set_processors_ _automatic_test_generation_ sinking 8-way hand-optimized _qwerty_-like _voltage-scaling_ _systematic_design_ performance-complexity	Cluster 3 (description: Computer Architecture) _fpga_devices_ _processor_arrays_ dma 8-node x-tree _heterogeneous_multiprocessors_ _back_end_ _fir_ spe openrisic
Cluster 2 (description:) mdst x_i _aggregate_signatures_ push-out _job_shop_ _forger_attack_ halts ssat prespecified _joint_source_	Cluster 4 (description:) _obstacle-avoiding_ nlog2 parentheses _outerplanar_graph_ _algebraic_number_field_ lim out-branching _expected_distortion_ hyperheuristic parameter	Cluster 6 (description:) _large_spatial_ _backscatter_ malignancy _gene_products_ pumps _expression_patterns_ _rna_secondary_structure_ per-electromyographic acute
Cluster 4 (description:) auctioneer/buyer _eye_tracking_ humans collected interview communicate _physical_objects_ _private_information_ weeks uncover	Cluster 7 (description:) _sdb_ kooperative ddc eng 5- fréquents dissections superstrings linearen _asp	Cluster 9 (description:) rserpool _dynamic_service_ _mobile-agent_based_ _management_issues_ pacts ome _educational_hypermedia_ inter-organization http-based database-driven

YELP.100K		
k= 5	k = 10	k = 25
Cluster 0 (description:) deterrent _great_pho_ exemplary _neapolitan_ _loose_leaf_ _previously_mentioned_ vastly _aloha_kitchen_ hh _pastys_	Cluster 0 (description: Time) _returned_home_ _tasty_meal_ extending _valentines_day_ _mon_-_thurs_ _friday_ tuesday luci's _friday_morning_ _seahawks_	Cluster 0 (description: Time) _friday_afternoon_ _sunday_nights_ _black_friday_ _patrick's_day_ _tuesday_morning_ _friends_birthday_ _mothers_day_ _busy_nights_ _weekend_night_ _lunch_hours_
Cluster 1 (description:) corroded leaky _classical_music_ _particle_board_ _construction_ blare _picnic_bench_ jugs _over-run_ _tobacco_	Cluster 2 (description:) _great_breakfast_ well- 83 _ipa's_ _foo-foo_ _organic_produce_ depending service- solid mocha's	Cluster 3 (description:) _middle_aged_ cliches woos relatives _nhl_ _grandparents_ spectators watch sex _business_travelers_
Cluster 3 (description: Feelings) _valentines_day_ hears unenthusiastic affiliation _handgun_ 5-minute sharpening _flaunt_ _low-cost_ growlers	Cluster 6 (description: Foods) _beer_glass_ _chocolate_raspberry_ cinnamony _navy_ rectangular avacados tangy _tender_chicken_ _red_wine_ icey	Cluster 2 (description:) hick knowingly optimist naw _woo_ pagers misread accomplished stephy fryers
Cluster 4 (description: Foods) _candy_cane_ salad _jollof_rice_ appys tsatziki rehydrated quirkiness jimmies greenbeans mix-n-_match_	Cluster 9 (description: Menus) trio baked wives dol _sizzling_rice_soup_ _holy_crap_ _marinated_beef_ _pulled_chicken_ _wontons_ _croissant_sandwich_	Cluster 6 (description:) domed embossed _stainless_steel_ fluorescent _large_bar_ equipped beaches courtyard _mini_bar_ shower