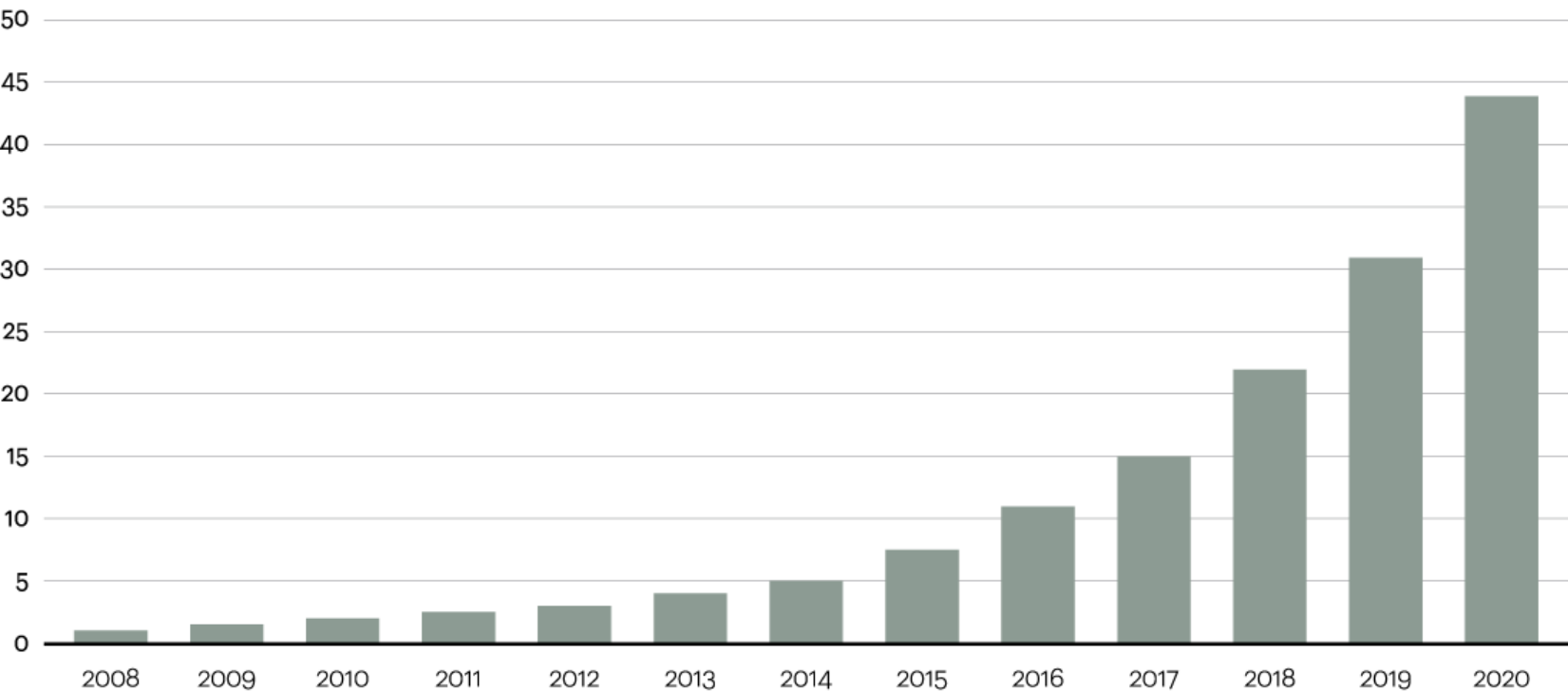


Distributed Representation of Documents with Explicit Explanatory Features

September 14th, 2015
SNU Data Mining Center
Han Kyul Kim

Figure 1
Data is growing at a 40 percent compound annual rate, reaching nearly 45 ZB by 2020

Data in zettabytes (ZB)

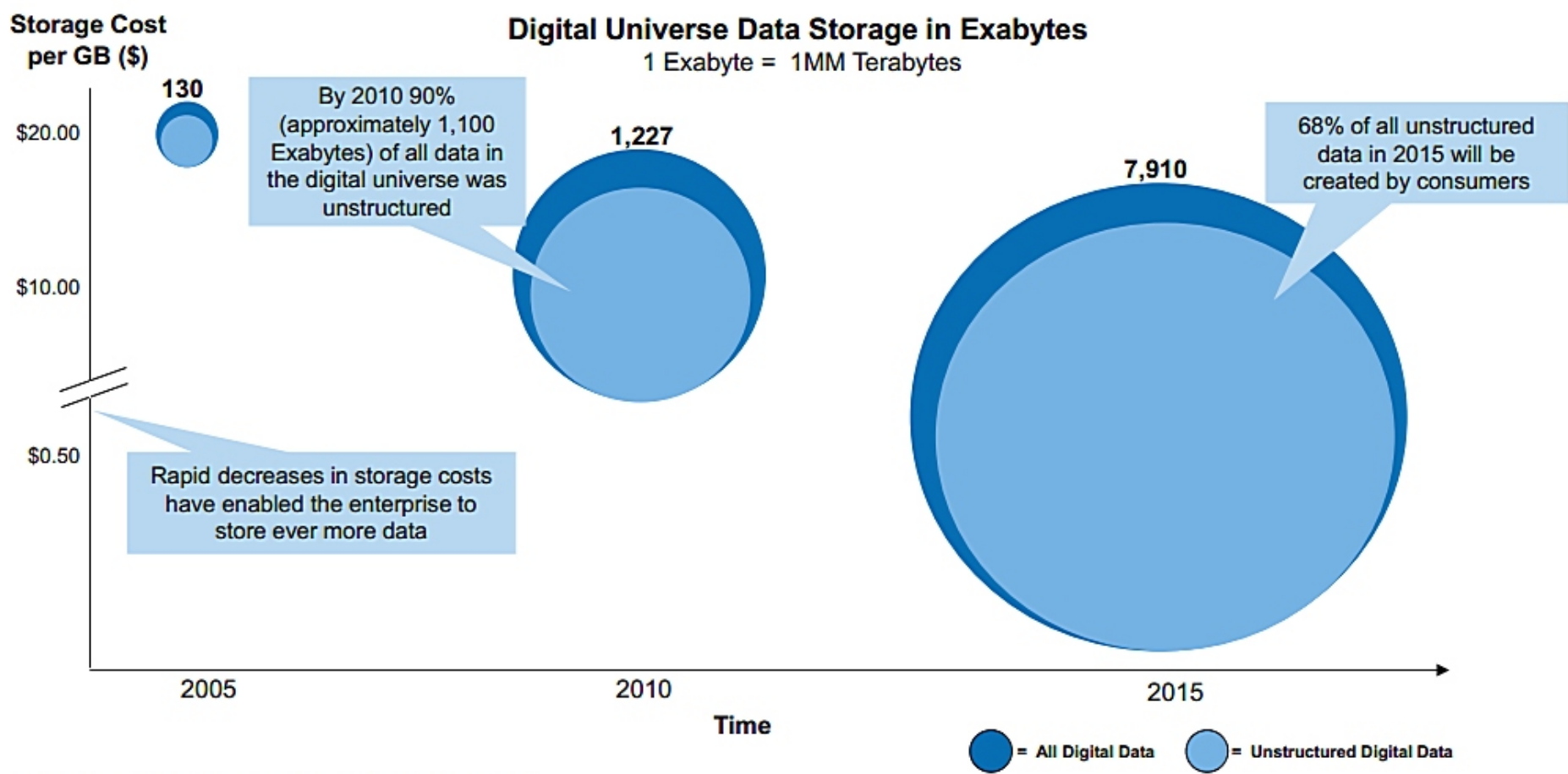


Source: Oracle, 2012

*** 1 ZB = 1 billion TB**

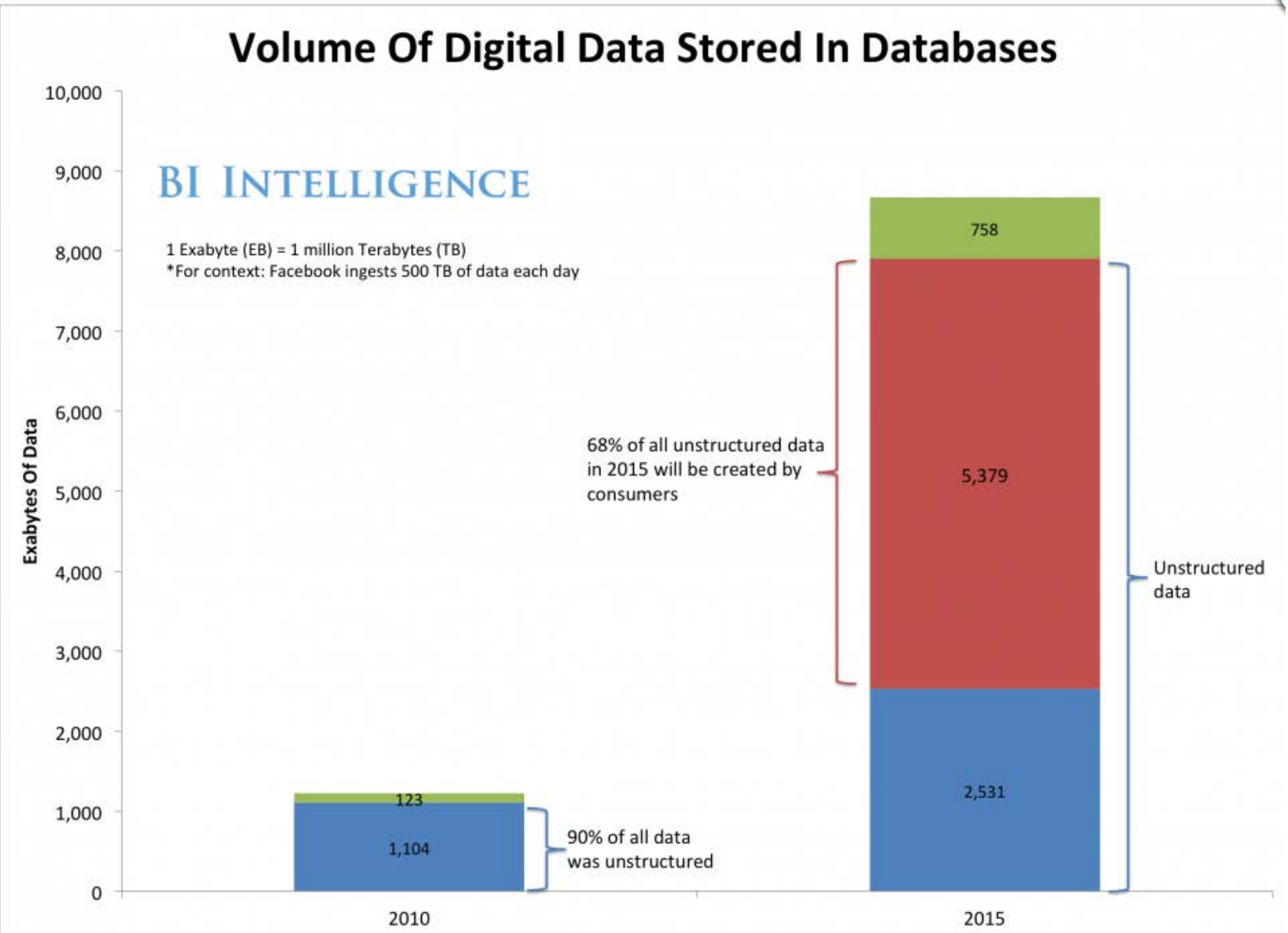
Ascendance of “Unstructured Data”

Unstructured consumer data, called Big Data, represents majority of growth in data volume, up 56% CAGR since 2005



Source: IDC's Digital Universe Study, sponsored by EMC, June 2011

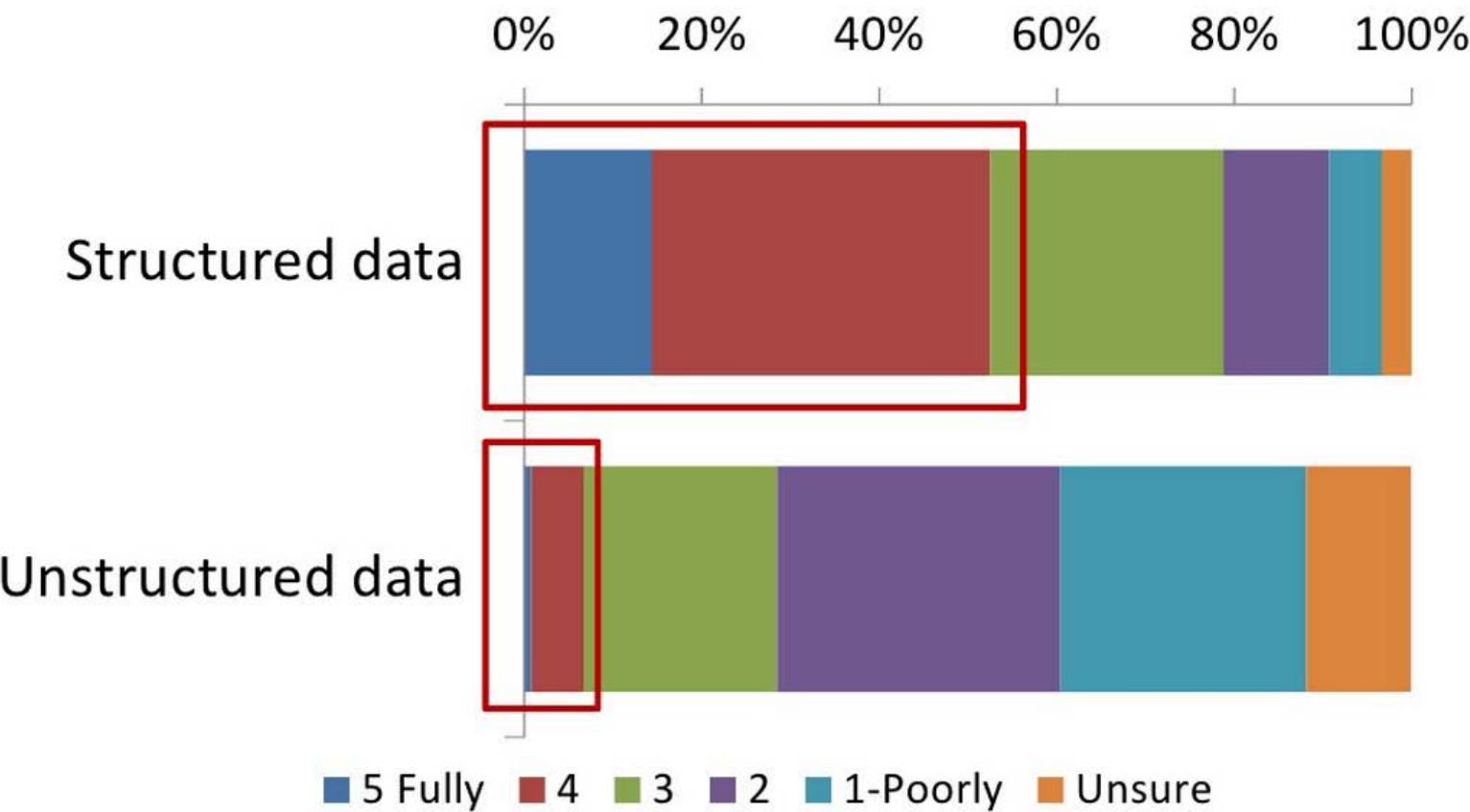
Ascendance of “Unstructured Data”



Source: IDC, BI Intelligence Estimates

Are companies ready?

Considered overall, to what degree does your organization exploit its information assets for analysis and decision making purposes?

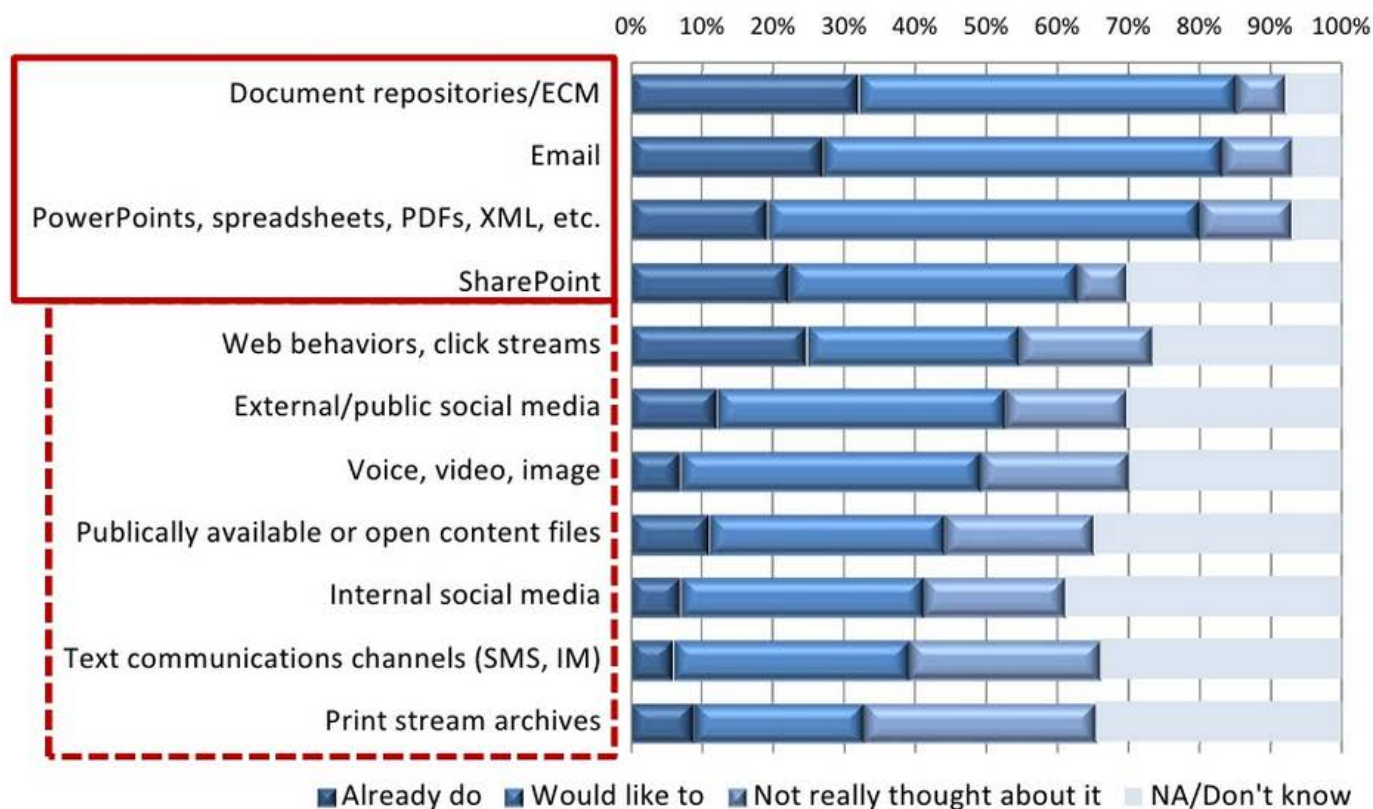


Source: Online survey of Register readers, 122 respondents, first half of November 2011, Freeform Dynamics



What kind of unstructured data are they interested in?

Are there large unstructured or semi-structured data repositories (ie, text, rich media, etc.) in your business that you would like to analyze, monitor or query - as opposed to search/retrieve?



Most popular are the basic ones – maybe some crossover “search/analyze.”

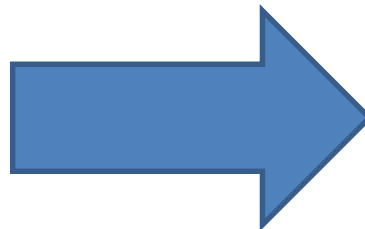
But 30-50% interested in a wide range of other applications.

But how should we represent a document?



- Number of Raw Words
- Number of Paragraphs
- Semantic Frequency
- Word Order
- Part of Speech
- Topic
- Sentiment

.....



Doc 1: [1, 4, 3, 0, ...]

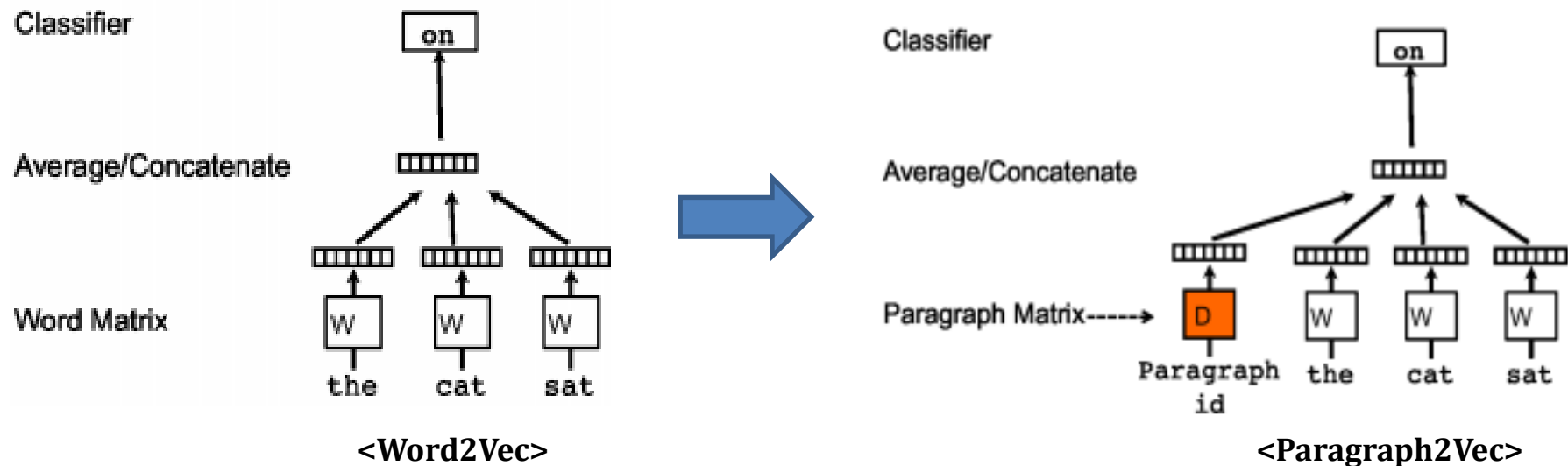
Doc 2: [5, 2, 8 0, ...]

.....

Term-Document Matrix Representation

- Bag of Word Hypothesis:
 - Frequencies of words in a document tend to indicate the relevance of documents
- Examples:
 1. Bag-of-words
 - Word order & semantic information are lost
 - Different sentences have exactly same representation as long as the same words are used
 2. Bag-of-n-grams
 - Restricted to short interval of texts
 - High dimensionality
 - Semantic information lost
 3. Weighted Average of word vectors
 - Loses the word order
 - Not that much different from BOW
 4. Word vectors + Parse Tree
 - Preserves the word order
 - Restricted to representing at the sentence level
- Yet, it still is one of the most popular method for representing document
 - Lucene: Term-Document Matrix method

Distributed Representation - Doc2Vec



Le, Q. V., & Mikolov, T. (2014). Distributed representations of sentences and documents

- Extension of Word2Vec
- Based on Distributional Hypothesis (Harris, 1954)
 - Words that occur in similar contexts tend to have similar meanings
- Able to represent input sequence of **variable** length
- Captures **semantic information** of words within paragraphs
- Doesn't rely on weight functions or parse trees
- Trained from **unlabeled data**

Next Generation of Document Representation?

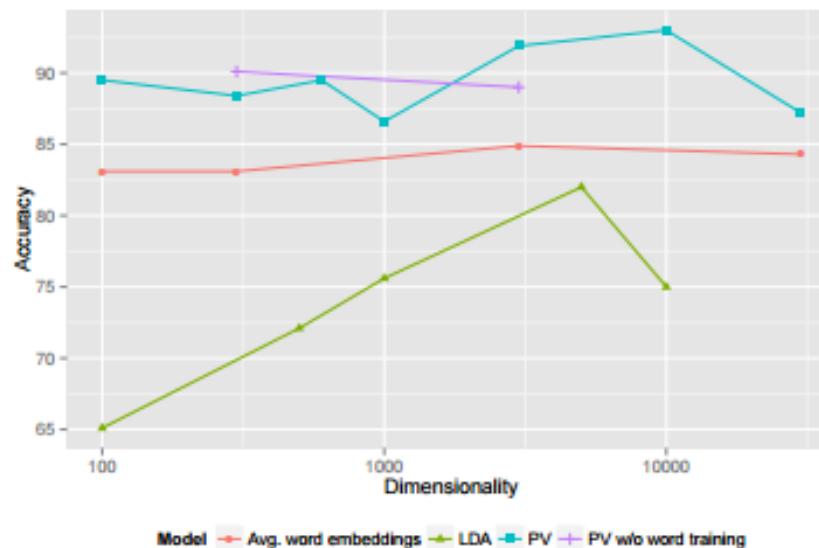


Figure 4: Results of experiments on the hand-built Wikipedia triplet dataset.

Table 3: Performances of different methods on hand-built triplets of Wikipedia articles on the best performing dimensionality.

Model	Embedding dimensions/topics	Accuracy
Paragraph vectors	10000	93.0%
LDA	5000	82%
Averaged word embeddings	3000	84.9%
Bag of words		86.0%

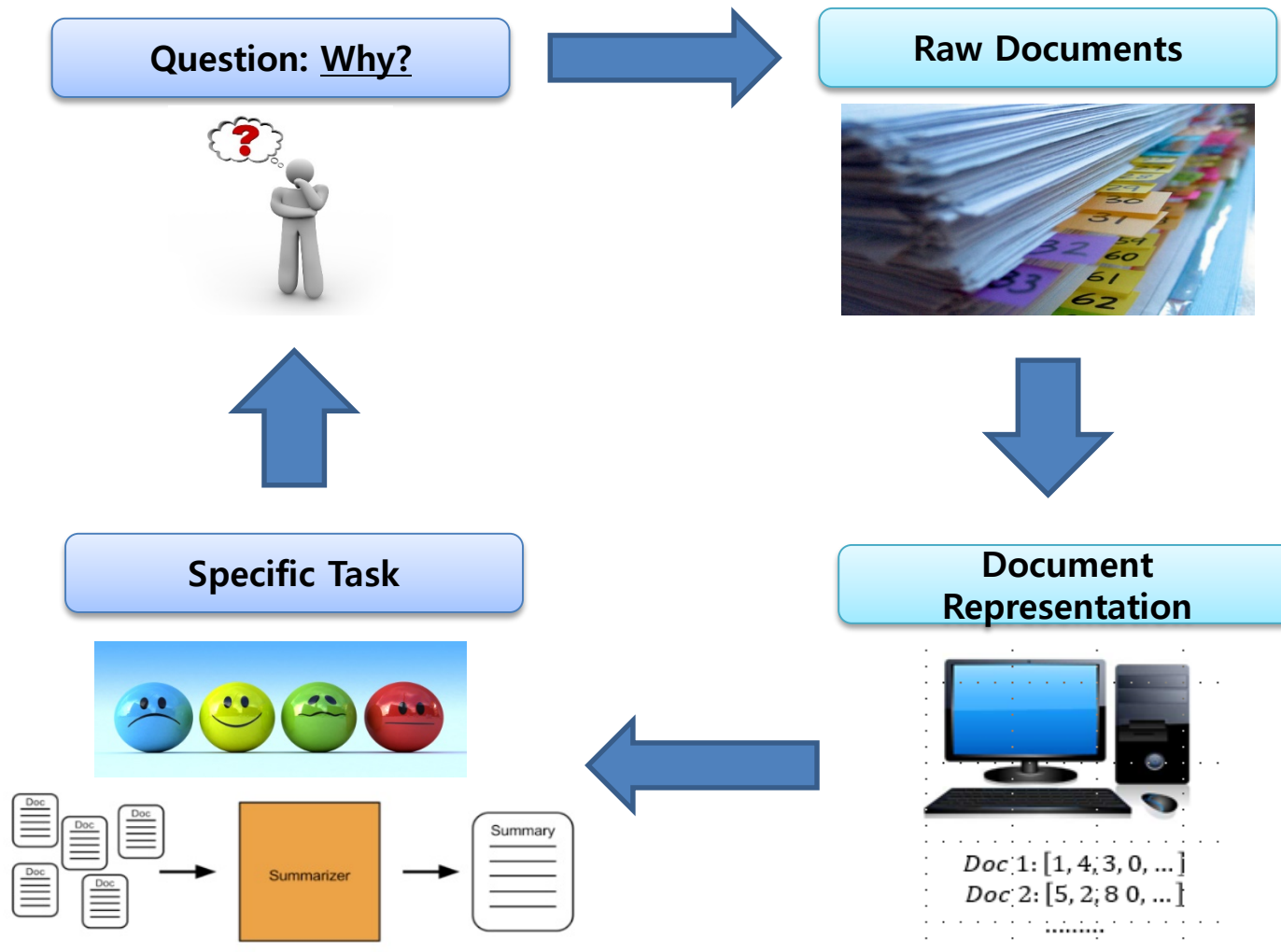
Dai, Andrew M., et al. "Document Embedding with Paragraph Vectors." NIPS Deep Learning Workshop. 2014.

Drawbacks of Doc2Vec

[-0.08759557455778122, -0.04312118515372276, -0.08494572341442108, -0.024585919454693794, -0.05785191431641579, -0.02659076638519764, 0.04704275727272034, -0.03940117731690407, 0.005195754114538431, -0.018994472920894623, -0.030896589159965515, -0.02599106915295124, -0.029802896882401276, -0.009517285041511059, -0.0362452491178818, 0.0029738633893430233, -0.04270448908209801, -0.0809769511461258, -0.04064304754137993, 0.017775749787688255, 0.0910411849617958, 0.05333533510565758, -0.07692492008209229, 0.08628936856985092, -0.042326122522354126, -0.007681592833250761, 0.0414172001183033, -0.030358949676156044, 0.05717118829488754, 0.0396726056933403, -0.09482961862945557, -0.05382954701781273, -0.016189705580472946, 0.001355066894526482, 0.004251557867974043, -0.10439810156822205, -0.01734139770269394, 0.08733568304539932, -0.02014184184372425, 0.06905293464660645, -0.052193693816661835, -0.008379205130040646, 0.050789929926395416, -0.0521097406744957, 0.02524719017673832, -0.09064795076847076, -0.01605154387652874, -0.08548879623413086, 0.09579522907733917, 0.07222563773393631, -0.01747663877904415, -0.07119490951299667, 0.04312814772129059, 0.006512957159429789, 0.04662078991532326, 0.053695641458034515, 0.0017072523135031462, 0.0468018501996994, 0.03436211124062538, -0.05560864135622978, 0.04996718838810921, 0.09512156248092651, 0.004873079247772694, -0.001316450640559196, 0.063605397939682, 0.02428655236364246, 0.009762034751474857, -0.04571126028895378, 0.032288745045661926, -0.06711595505475998, 0.03179262951016426, 0.00480815464087948203, 0.013781636022031307, 0.049716588109731674, -0.04940832778811455, 0.04659050330519676, 0.08707352727651596, -0.10198234766721725, 0.0012964112684130669, 0.019826047122478485, 0.020798761397600174, 0.0398060940260509, -0.01610562391579151, 0.09880314767360687, 0.035302355885505676, -0.03354038670659065, -0.060332611203193665, 0.009992017410695553, -0.07922962307929993, -0.04672875627875328, 0.02949077662825584, 0.007386107929050922, 0.014454830437898636, 0.03679213672876358, 0.02259526588022709, -0.07544633001089096, 0.03752242028713226, -0.04298776760697365, 0.04293576627969742, -0.04286639392375946, 0.05493065342307091, -0.010759598997368813, -0.02642848715186119, 0.05936839431524277, 0.01073556483268738, -0.022106561809778214, 0.005098363384604454, 0.02817825973033905, 0.06781460344791412, 0.01152470614771141, 0.04529837146401405, -0.10414177924394608, -0.06333499401807785, -0.025369135662913322, 0.0305448077619076, 0.08876033872365952, -0.03237186372280121, 0.08923118561506271, 0.035788971930742264, -0.0735849142074585, -0.11431705206632614, 0.004165078978985548, 0.03989437595009804, 0.011698578484356403, 0.009127369150519371, -0.007194945588707924, -0.06279811263084441, -0.012036782689392567, 0.040613592482089996, -0.07241667807102203, 0.06803597509860992, 0.03180589095059395, 0.02011158317274994, 0.05419156327843666, -0.06192755699157715, 0.03006706014275551, -0.015478908084332943, -0.0518170624109703, -0.027904511498046875, 0.01615394838154316, 0.0011802234221249819, -0.052437208592689169, -0.058099536028824, 0.025693533942103386, -0.04456381872296333, -0.05235142633318901, 0.002849023789167404, 0.02517843246459961, 0.07427450269460678, -0.0041514914482831955, 0.06708431243896484, 0.010245034471154213, -0.037210747599601746, -0.05194145068526268, 0.011150983162224293, -0.01369072962552309, -0.004474998451769352, 0.08407635986804962, 0.01842060498893261, 0.06935162097215652, -0.070859894156456, -0.04950730875134468, -0.024812163114548, -0.006143994629383087, 0.018353456631302834, 0.10588039457798004, 0.048505134880542755, -0.014145134948194027, -0.02343590557575226, -0.0077395206317305565, -0.024359061208343506, 0.02147112786769867, 0.07666371762752533, 0.06256308406591415, 0.012278590351343155, -0.017493827268481255, -0.025852585211396217, 0.02850286103785038, -0.04092925786972046, 0.1389494389295578, -0.0005276211304590106, 0.019603941589593887, 0.047872498631477356, -0.02479643188416958, -0.03278869017958641, -0.03298942744731903, -0.0749734565615654, 0.017869295552732932, 0.014359974302351475, 0.03260161727666855, -0.05898589268326759, 0.049986109137535095, -0.016536226496100426, -0.019838936362691163, -0.030345618724822998, -0.046361502259896971, -0.07739298790693283, 0.05453772097826004, -0.015122903510928154, -0.055823419243097305, 0.02098872688843479, 0.0380556583404541, 0.0410100519657135, 0.06500554829835892, -0.014411268755793571, 0.06324303895235062, 0.12235700339078003, -0.003138007363304496, 0.04150061309337616, -0.05469806492328644, 0.028155071660876274, 0.02540171346084118, 0.02362082526087761, -0.055749375373125076, 0.034880928695201874, 0.029251024172700653, -0.044924408197402954, -0.05796622857451439, 0.16091801226139069, 0.013727911747992039, -0.002823801381855297, 0.0050589111633598804, 0.009330375120043755, -0.040790800005197525, 0.012631899677217007, -0.03498981520533562, -0.09075972437858582, -0.04068059101700783, 0.056468620896339417, 0.0793653130531311, -0.10637512058019638, 0.024917149916291237, 0.076641976833435, 0.0889914280939102, -0.0765515521764755, 0.018688995391130447, -0.03665241599082947, -0.007174948696047068, -0.026764938607811928, 0.006813205778598785, 0.03434935212135315, 0.06390520185232162, 0.005475071258842945, 0.006454255897551775, 0.01241873949766159, -0.04875058806727814, -0.025715982541441917, -0.0013290401548147202, -0.0036538743879646063, -0.03493126481771469, 0.07465724647045135, 0.04718988761305809, -0.027499066665768623, 0.011664669029414654, 0.020739786326885223, -0.001199969439767301, 0.022283419854938984, -0.039574239403009415, 0.027619406580924988, 0.01381973270368528, 0.00963809332549572, 0.017708426634141445, -0.013031989336013794, -0.0582968100905418, -0.002542087387293577, 0.03596331179141998, 0.0165165513753891, -0.031340003236148026, -0.03133828236165047, 0.07702472805976868, -0.06405989825725555, -0.05113033577799797, 0.11122194677591324, -0.0431693010032177, 0.0351727195084095, -0.06082572788000107, -0.05721887946128845, -0.04940216988325119, -0.04693160206079483, -0.04560410603808802, -0.06396282464265823, -0.08668574690818787, -0.008351934142410755, 0.08004070073366165, 0.020836271345615387, 0.03663250431418419, 0.01230071671370807, -0.042455870658159256, 0.0060779107204937935, 0.0036929140333086252, 0.004990659188479185, -0.00278487685136497, -0.03410743291995163, -0.058781374245882034, 0.050218498492240906, 0.018009591847658157, 0.05601579323410988, 0.06174979358911514, 0.009586947038769722, 0.10309688746929169, -0.0013724834425374866, -0.037349410355091095, -0.06568935513496399, -0.0019719060510396957, -0.02522912435233593, -0.006819620728492737, -0.08099189400672913, 0.14139385521411896, 0.02598827298259735, 0.044450853019952774, 0.09344274550676346, -0.022272149108999696, -0.019451434918618202, 0.055291395634412766, -0.03372661769390106, 0.07417678833007812, 0.0070419879630208015, -0.053876567631959915, 0.10055916011324366, -0.11293848603963852, -0.03376157209277153, 0.060266170650720596, 0.0392690044117464, -0.035960644483566284, 0.008666090667247772, 0.027477947995066643, 0.03710789978504181, -0.023510152474045753, 0.03532886877655983, -0.002786519005894661, -0.021855372935533524, -0.028223032131791115, 0.042905956506729126, -0.02157396450638771, -0.02910805597127676, 0.039419323205947876, -0.06341256201767242, -0.04983802139759064, 0.04110337421298027, -0.08459826558828354, -0.0303527776812315, 0.03796708953600817, 0.17562396824359804, 0.05129672959446907, 0.005838119424879551, -0.05645650997757912, 0.02716675214469433, -0.024479210376739502, -0.05037922411257059, 0.04924848675727844, 0.016620146110653877, -0.04153639078140259, 0.02921335957944393, 0.040044404566287994, 0.012756554409861565, -0.007722984999418259, -0.0386621430516243, 0.017636995762586594, -0.04599433019757271, -0.05585790053009887, -0.00897784624248743, -0.07755934474050095, 0.02797515639378071, 0.05918818712234497, 0.00961016118526487, 0.09371864050626755, 0.006441071629524231, -0.015960507094860077, 0.10366759449243546, 0.00836758129298687, 0.05804479494690805, 0.04077683761715889, -0.003010124666613206, -0.007653041727080774, 0.12586210668087006, 0.02910740114748478, 0.04234848675727844, 0.01783227137386818, 0.07243572175502777, 0.03400840610265732, 0.027588292956352234, -0.0224411990493536, -0.029101211577653885, 0.00895093847066164, -0.013047640211880202, 0.05406863187909126, -0.001707519079002738, -0.0310812415007086754, -0.06595993787050247, 0.01366981189174652, -0.019473319873213768, -0.07515253126621246, 0.05773814767599106, 0.005721567664295435, 0.02362120896577835, 0.040131740272045135, 0.07627321779727936, -0.04665115848183632, -0.043376557528972626, 0.061202943325042725, 0.004489867947959563, -0.020421821624040604, -0.044818829745054245, -0.0038449352141469717, -0.05663484334945679, 0.03718886151909828, -0.02091551127433777, -0.03713458776473999, -0.05943121761083603, -0.06698095798492432, 0.012730694375932217, -0.021541433408856392, -0.018205733969807625, 0.014323082752525806, -0.07051131874322891, 0.11712339520454407, 0.000716670940644766, -0.06536946445703506, -0.015314917096813679, 0.034568077110634, 0.061441387981176376, 0.024788793176412582, 0.01627667434513569, -0.06662845611572266, -0.015765206888318062, -0.0037493400741368532, -0.004942044615745544, -0.042015254497528076, -0.07548461109399796, -0.03183511644601822, 0.061288982629776, -0.06246870756149292, -0.0013078112388029695, -0.06325135380029678, -0.12265635821801224, 0.020725131034851074, 0.028045836836099625, 0.0168102215975523, -0.014592664316295677, 0.043280523270368576, -0.

Having a good representation form itself is not **the ultimate goal** of document representation!

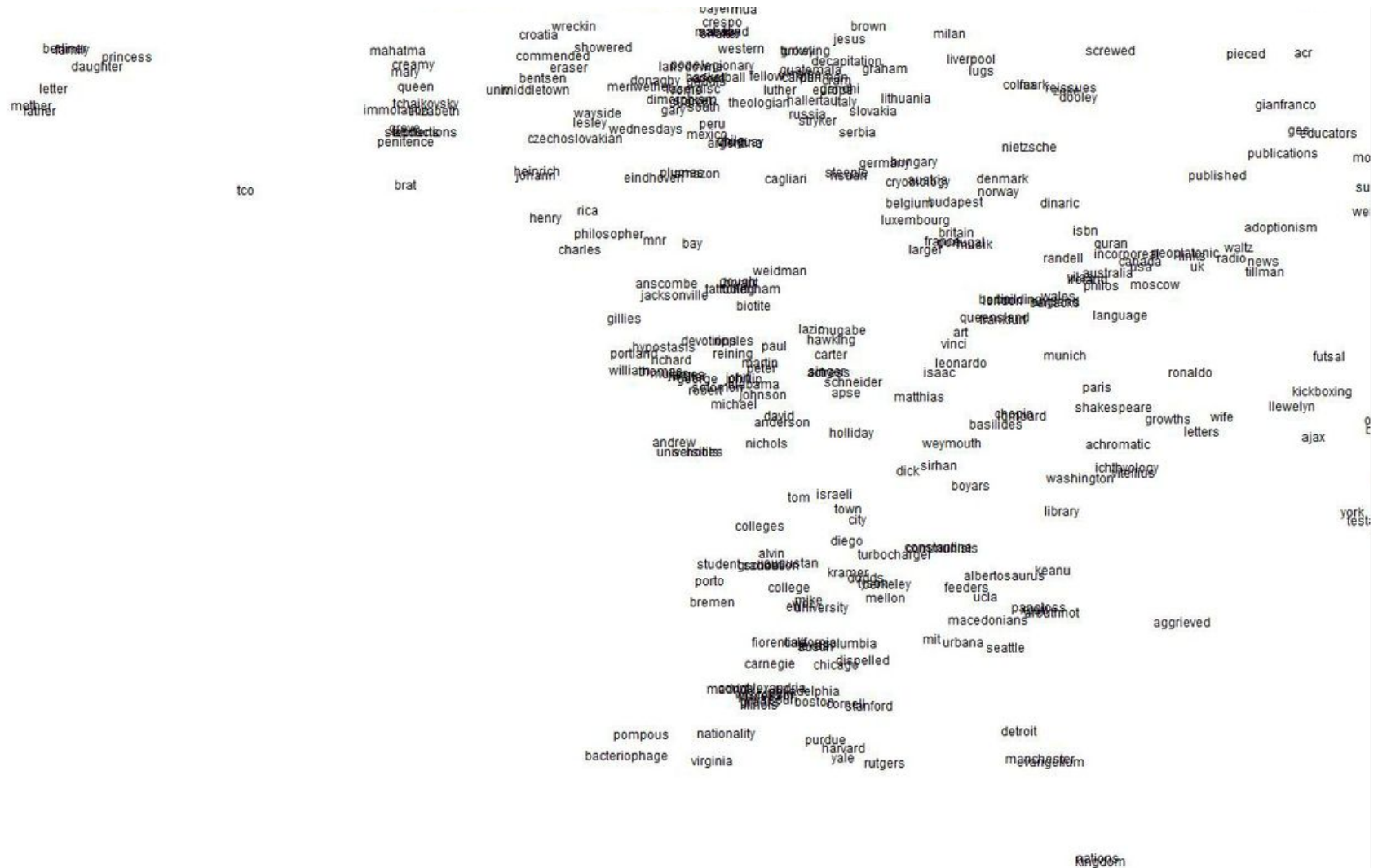
Drawbacks of Doc2Vec



Doc2Vec Representation provides with an insight to **what's** happening amongst the documents, but **not how and why** it's happening

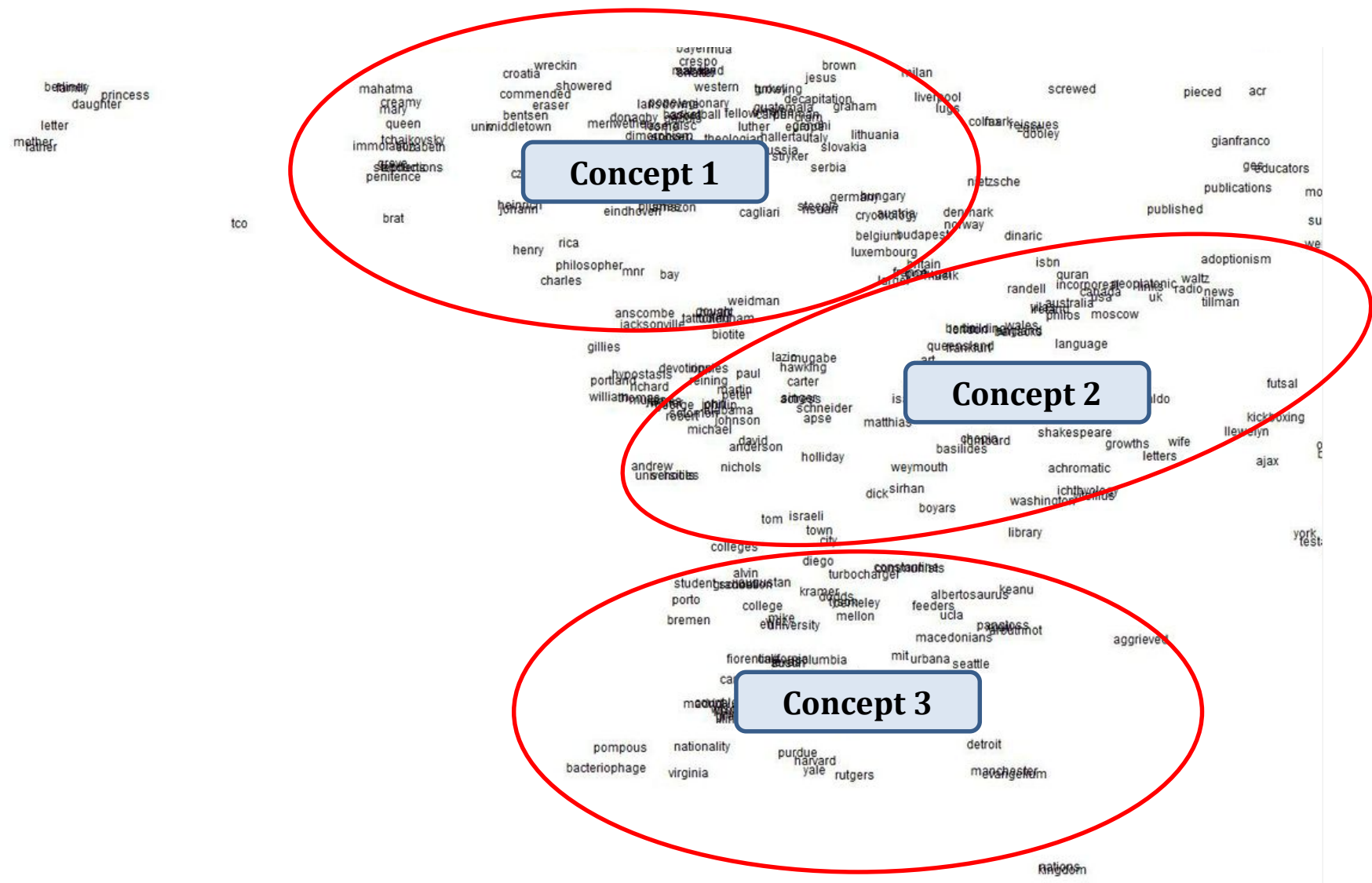
Proposed Framework

1. Train Word2Vec with the collection of documents

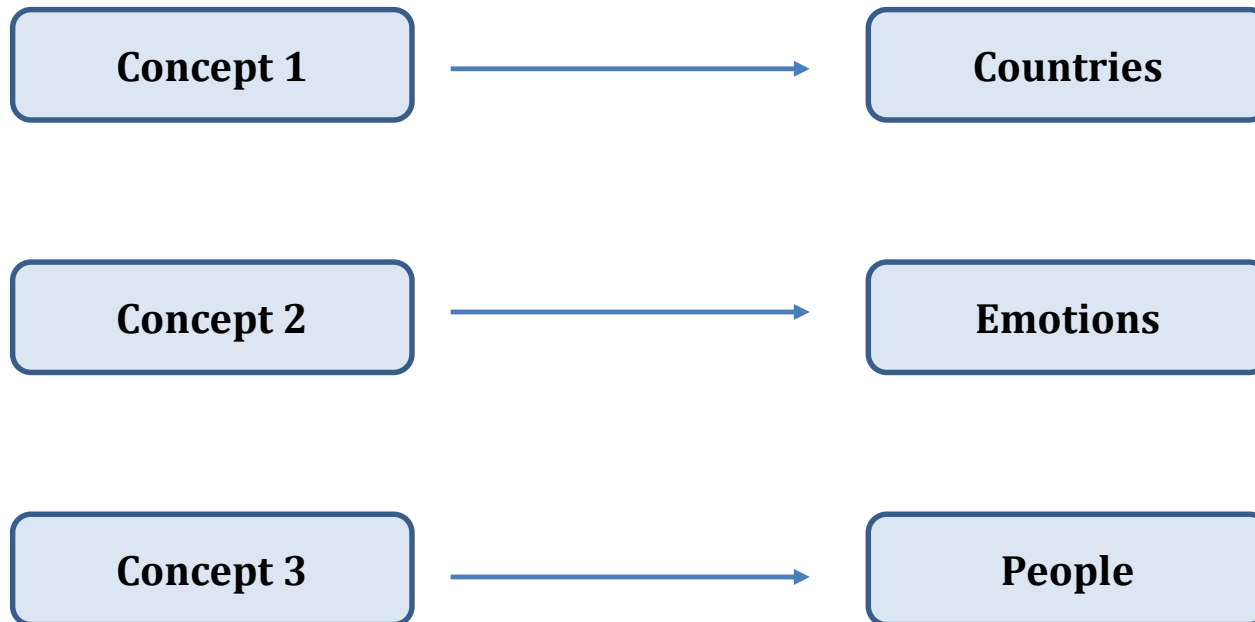


Proposed Framework

2. Cluster word2vec generated vectors to create clusters of concepts

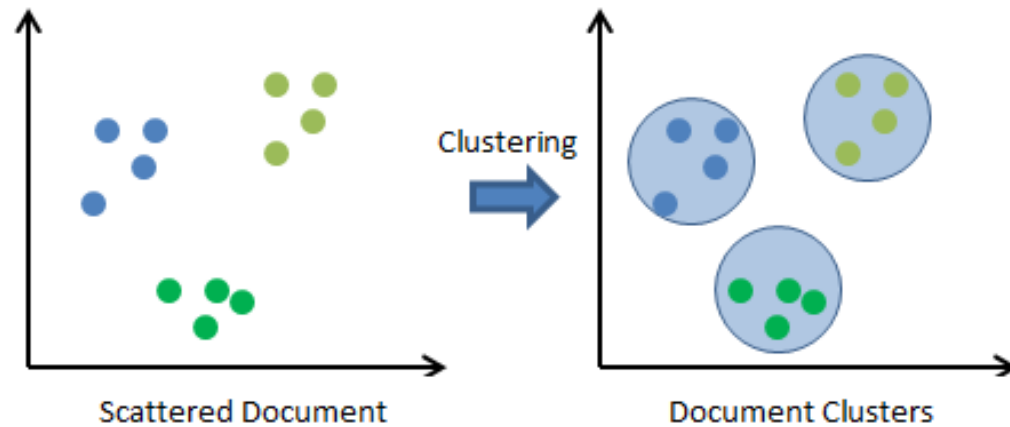


3. Label the concepts using the words associated with each cluster



Proposed Framework

5. Test the effectiveness of the document representation through document clustering and classification



This proposed framework...

- Maintains the representational power of the distributed representation, while providing explicit and intuitive features
- Provides reasons and logic behind the representation
- Includes more holistic semantic information about documents
- Suggests language independent framework for representing documents

Reference

Dai, Andrew M., et al. "Document Embedding with Paragraph Vectors." NIPS Deep Learning Workshop. 2014.

Le, Quoc V., and Tomas Mikolov. "Distributed representations of sentences and documents." arXiv preprint arXiv:1405.4053 (2014).

Mikolov, Tomas, et al. "Distributed representations of words and phrases and their compositionality." Advances in Neural Information Processing Systems. 2013.

Mikolov, Tomas, Quoc V. Le, and Ilya Sutskever. "Exploiting similarities among languages for machine translation." arXiv preprint arXiv:1309.4168(2013).

Rong, Xin. "word2vec Parameter Learning Explained." arXiv preprint arXiv:1411.2738 (2014).

Turney, Peter D., and Patrick Pantel. "From frequency to meaning: Vector space models of semantics." Journal of artificial intelligence research 37.1 (2010): 141-188.