
BLG 540E
TEXT RETRIEVAL SYSTEMS

Dictionaries and tolerant retrieval

Arzucan Özgür



Announcements





Paper Presentations

- ▶ Please email me (ozgura@itu.edu.tr) your top 3 paper choices for presentation. Email me in advance to check if the papers are appropriate.
 - ▶ Due date: **03/03/2011 (Thursday) – 17:00.**
- ▶ You can choose from the suggested papers. But **at least one** of your three preferred papers should be proposed by you.
- ▶ I will try to assign you one of your three choices.
- ▶ You will have **20 minutes** for presentation + **5 minutes** for questions/discussion.



Some paper suggestions

- ▶ The XML Retrieval chapter in the book.
- ▶ Dafna Shahaf and Carlos Guestrin. Connecting the dots between news articles. In KDD '10: Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 623–632, New York, NY, USA, 2010. ACM.
- ▶ Turning Down the Noise in the Blogosphere, KDD2009.
 - ▶ <http://www.cs.cmu.edu/~dshahaf/kdd2009-elarini-veda-shahaf-guestrin.pdf>
- ▶ Michele Banko and Oren Etzioni. The tradeoffs between open and traditional relation extraction. In Proceedings of ACL-08: HLT, pages 28–36, Columbus, Ohio, June 2008. Association for Computational Linguistics.
- ▶ - Science Maps: Chen, C. (2006). CiteSpace II: Detecting and visualizing emerging trends and transient patterns in scientific literature. Journal of the American Society for Information Science and Technology, 57(3), 359-377. doi: 10.1002/asi.20317
 - ▶ <http://cluster.cis.drexel.edu/~cchen/citespace/doc/jasist2006.pdf>



Some paper suggestions

- ▶ Monika Rauch Henzinger, Finding near-duplicate web pages: a large-scale evaluation of algorithms, SIGIR, 2006, pp. 284–291.
- ▶ Andrei Z. Broder, Identifying and filtering near-duplicate documents, CPM, 2000, pp. 1–10.
- ▶ Gunes Erkan and Dragomir R. Radev. Lexrank: Graph-based centrality as salience in text summarization. Journal of Artificial Intelligence Research (JAIR), 2004.
- ▶ Qiaozhu Mei, Dengyong Zhou, Kenneth Church. Query Suggestion Using Hitting Time, Proceedings of the 16th ACM International Conference on Information and Knowledge Management (CIKM'08), pages 469-478, 2008.
- ▶ Patterns of Cascading Behavior in Large Blog Graphs by J. Leskovec, M. McGlohon, C. Faloutsos, N. Glance, M. Hurst. SIAM SDM 2007.
- ▶ R.W.White and S. M. Drucker. Investigating behavioral variability in web search. In WWW'07.
- ▶ E. Agichtein, E. Brill, and S. Dumais. Improving web search ranking by incorporating user behavior information. In Proceedings of ACM SIGIR 2006.



Some paper suggestions

- ▶ G. Cong, L. Wang, C.-Y. Lin, Y.-I. Song, and Y. Sun. Finding question-answer pairs from online forums. In Proc. 31st SIGIR.
 - ▶ <http://research.microsoft.com/en-us/um/people/cyl/download/papers/SIGIR2008-Gao-MSRA.pdf>
- ▶ Chirita, P.-A., Firan, C. S., and Nejdl, W. Personalized query expansion for the web. In SIGIR (2007), pp. 7-14.
- ▶ Collins-Thompson, K., and Callan, J. Query expansion using random walk models. In CIKM (2005), pp. 704-711.
- ▶ Pavel Calodo et al. Combining link-based and text-based methods for web document classification. CIKM 2003.



Some paper suggestions

- ▶ Search advertising using Web relevance feedback. Broder, P. Ciccolo, M. Fontoura, E. Gabrilovich, V. Josifovski, and L. Riedel. (CIKM, 2008)
- ▶ Automatic Generation of Bid Phrases for Online Advertising. Ravi, S.; Broder, A.; Gabrilovich, E.; Josifovski, V.; Pandey, S.; Pang, B. WSDM (2010)
- ▶ Using Landing Pages for Sponsored Search Ad Selection. Choi, Y.; Fontoura, M.; Gabrilovich, E.; Josifovski, V.; Mediano, M.; Pang, B. (WWW 2010)
- ▶ Ganesh Ramakrishnan, Soumen Chakrabarti, Deepa Paranjpe, and Pushpak Bhattacharya. Is question answering an acquired skill? In Proceedings of the 13th international conference on World Wide Web, 2004.



Project

- ▶ Please submit a one-page project proposal by e-mail.
 - ▶ Due date: **09/03/2011 Wednesday 17:00.**
- ▶ You can choose from the list of project ideas or propose your own.
- ▶ You can work in teams of two, or individually.
- ▶ The last two weeks (06/05/2011 and 13/05/2011) will be allocated for your project presentations (15 min. presentation + 5 min. questions/discussion).



Some Project Ideas

- ▶ Build a question answering system.
- ▶ Build a language identification system.
- ▶ Social network analysis from text
- ▶ Query log analysis.
- ▶ information extraction
- ▶ information extraction for biology (e.g. extracting protein interactions)
- ▶ blog analysis
- ▶ Sentiment/polarity extraction
- ▶ document duplicate and near-duplicate recognition
- ▶ clustering scientific papers by topic using citation information



Some Project Ideas

- ▶ automatic query correction/expansion
- ▶ query completion/recommendation
- ▶ extract names of people and their descriptions from the web
- ▶ finding similar documents
- ▶ named entity recognition
- ▶ named entity disambiguation
- ▶ movie recommendations
- ▶ adversarial IR (spam)
- ▶ language modeling for IR
- ▶ Text classification/clustering
- ▶ summarization





Recap



Type/token distinction

- **Token** – an instance of a word or term occurring in a document
- **Type** – an equivalence class of tokens
- *In June, the dog likes to chase the cat in the barn.*
- 12 word tokens, 9 word types

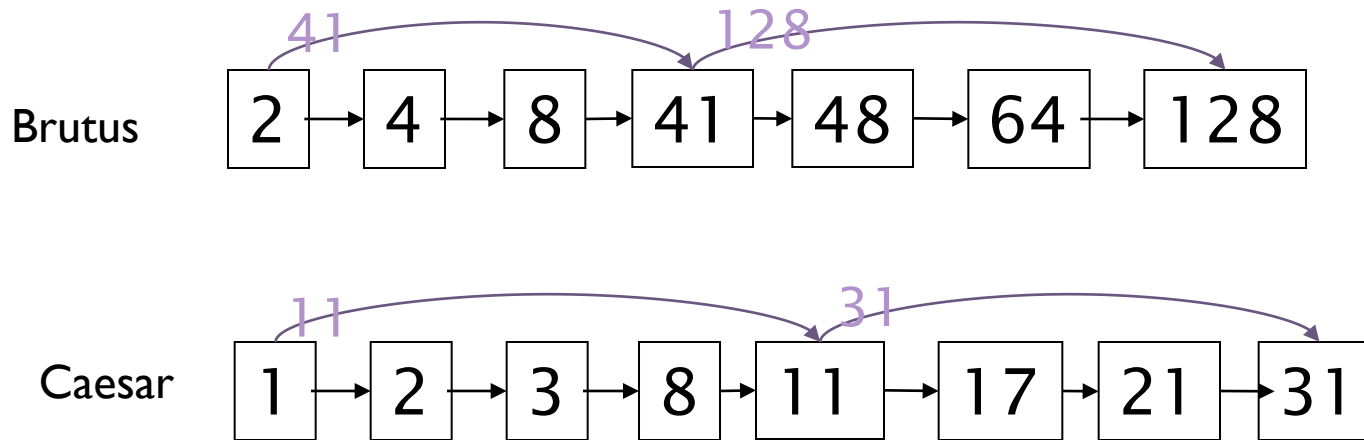
Problems in tokenization

- What are the delimiters? Space? Apostrophe? Hyphen?
- For each of these: sometimes they delimit, sometimes they don't.
- No whitespace in many languages! (e.g., Chinese)
- No whitespace in Dutch, German, Swedish compounds (*Lebensversicherungsgesellschaftsangestellter*)

Problems with equivalence classing

- A term is an equivalence class of tokens.
- How do we define equivalence classes?
- Numbers (3/20/91 vs. 20/3/91)
- Case folding
- Stemming, Porter stemmer
- Morphological analysis
- Equivalence classing problems in other languages
 - More complex morphology than in English
 - Finnish: a single verb may have 12,000 different forms
 - Accents, umlauts

Skip pointers



Positional indexes

- Postings lists in a **nonpositional** index: each posting is just a docID
- Postings lists in a **positional** index: each posting is a docID and **a list of positions**
- Example query: “ to_1 be_2 or_3 not_4 to_5 be_6 ”

TO, 993427:

1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
2: $\langle 1, 17, 74, 222, 255 \rangle$;
4: $\langle 8, 16, 190, 429, 433 \rangle$;
5: $\langle 363, 367 \rangle$;
7: $\langle 13, 23, 191 \rangle$; . . .

BE, 178239:

1: $\langle 17, 25 \rangle$;
4: $\langle 17, 191, 291, 430, 434 \rangle$;
5: $\langle 14, 19, 101 \rangle$; . . . Document 4 is a match!

Positional indexes

- With a positional index, we can answer **phrase queries**.
- With a positional index, we can answer **proximity queries**.

Today's Lecture

- **Tolerant retrieval**: What to do if there is no exact match between query term and document term
- Wildcard queries
- Spelling correction





Dictionaries



Inverted index

For each term t , we store a list of all documents that contain t .

BRUTUS →

1	2	4	11	31	45	173	174
---	---	---	----	----	----	-----	-----

CAESAR →

1	2	4	5	6	16	57	132	...
---	---	---	---	---	----	----	-----	-----

CALPURNIA →

2	31	54	101
---	----	----	-----

⋮


dictionary


postings

Dictionaries

- The dictionary is the data structure for storing the term vocabulary.
- **Term vocabulary**: the **data**
- **Dictionary**: the **data structure** for storing the term vocabulary

A naïve dictionary

- ▶ An array of struct:

term	document frequency	pointer to postings list
a	656,265	→
aachen	65	→
...
zulu	221	→

char[20]

20 bytes

int

4/8 bytes

Postings *

4/8 bytes

- ▶ How do we store a dictionary in memory efficiently?
- ▶ How do we quickly look up elements at query time?



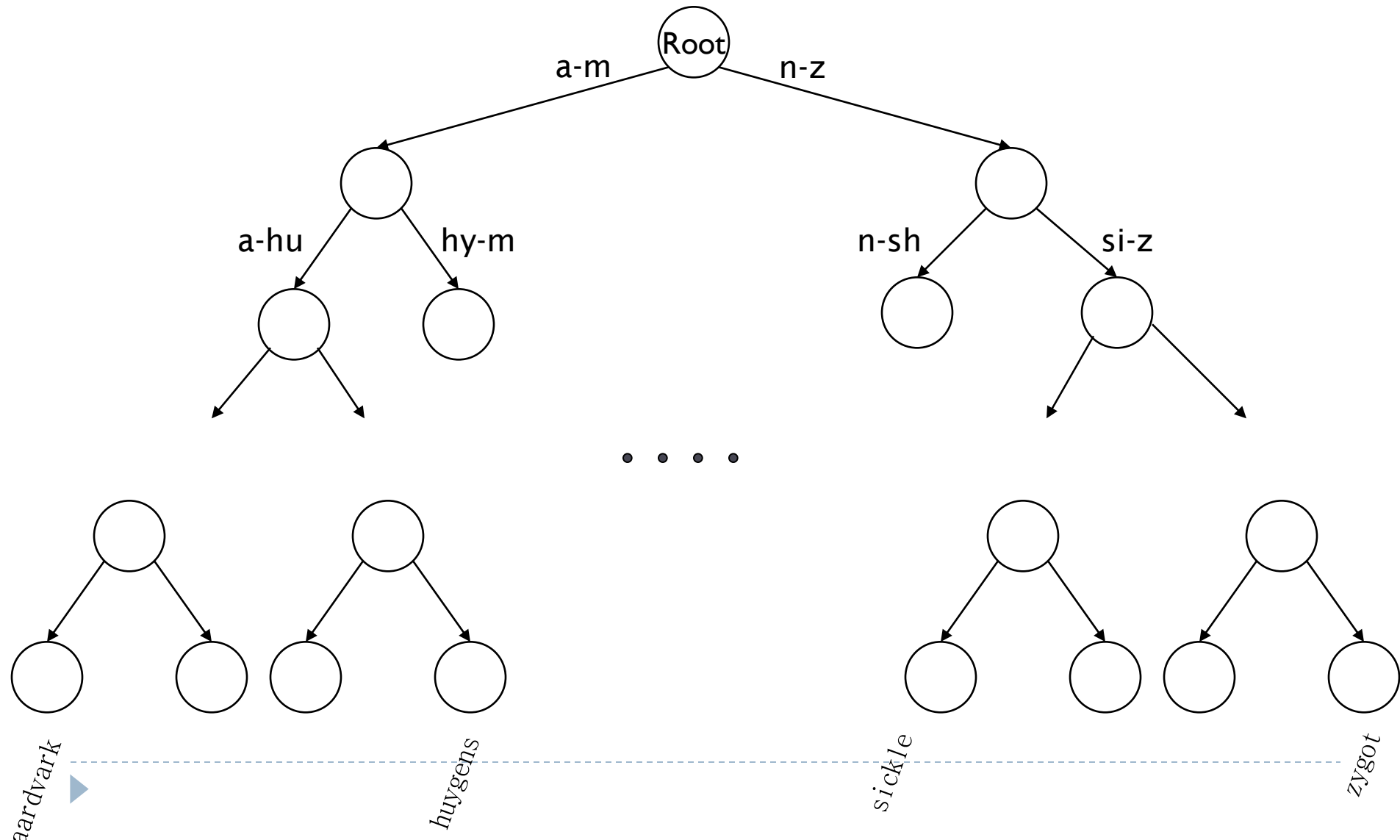
Data structures for looking up term

- Two main classes of data structures: hashes and trees
- Some IR systems use hashes, some use trees.
- Criteria for when to use hashes vs. trees:
 - Is there a fixed number of terms or will it keep growing?
 - What are the relative frequencies with which various keys will be accessed?
 - How many terms are we likely to have?

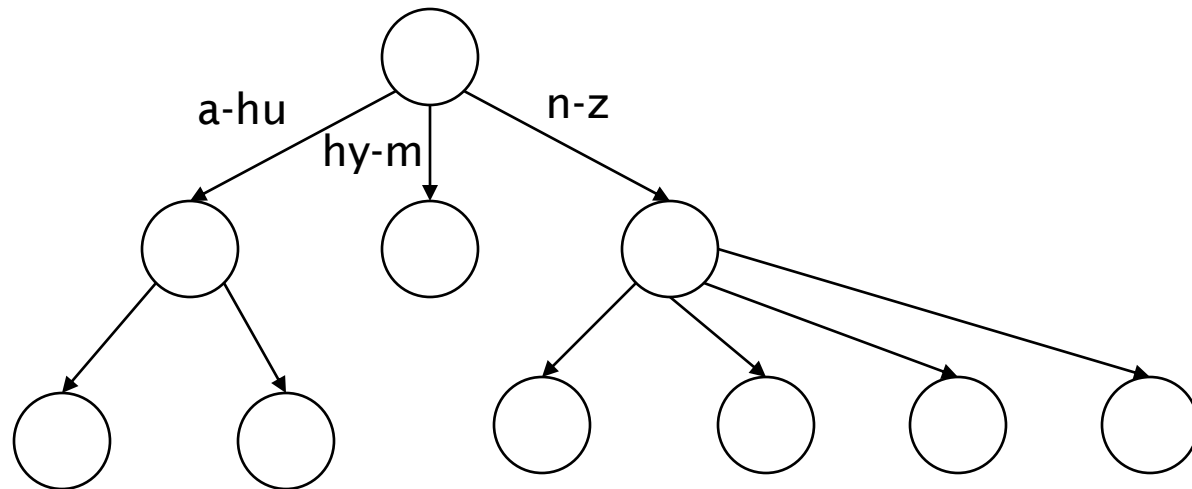
Hashes

- Each vocabulary term is hashed into an integer.
- Try to avoid collisions
- At query time, do the following: hash query term, resolve collisions, locate entry in fixed-width array
- Pros: Lookup in a hash is faster than lookup in a tree.
 - Lookup time is constant.
- Cons
 - no way to find minor variants (*resume* vs. *résumé*)
 - no prefix search (all terms starting with *automat*)
 - need to rehash everything periodically if vocabulary keeps growing

Tree: binary tree



Tree: B-tree



- Definition: Every internal node has a number of children in the interval $[a,b]$ where a, b are appropriate natural numbers, e.g., $[2,4]$.
-

Trees

- ▶ Simplest: binary tree
- ▶ More usual: B-trees
- ▶ Trees require a standard ordering of characters and hence strings ... but we standardly have one
- ▶ Pros:
 - ▶ Solves the prefix problem (terms starting with *hyp*)
- ▶ Cons:
 - ▶ Slower: $O(\log M)$ [and this requires *balanced* tree]
 - ▶ Rebalancing binary trees is expensive
 - ▶ But B-trees mitigate the rebalancing problem



Wild-card queries





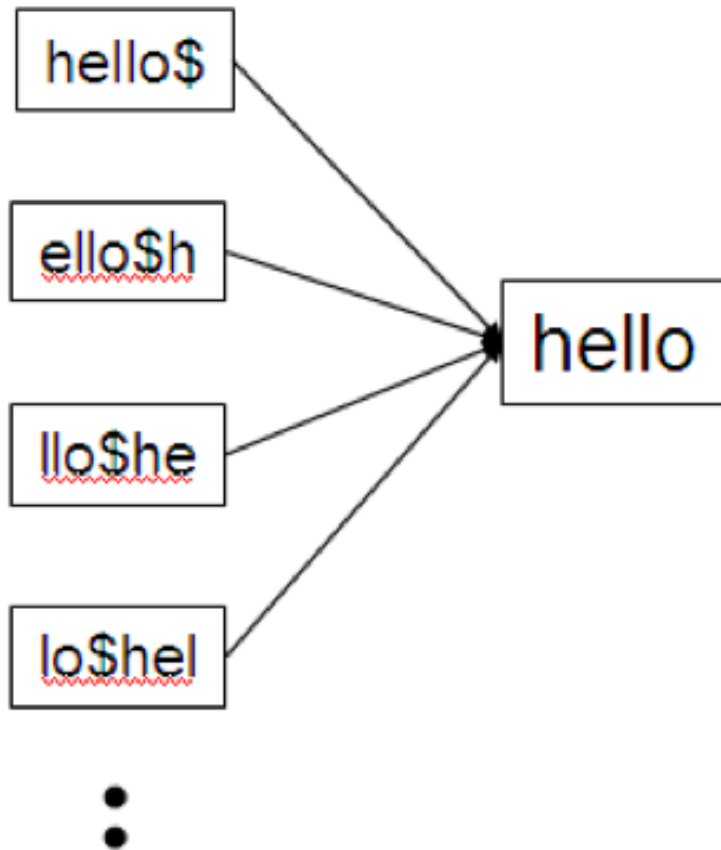
Wildcard queries

- mon^* : find all docs containing any term beginning with *mon*
- Easy with B-tree dictionary: retrieve all terms t in the range: $\text{mon} \leq t < \text{moo}$
- $^*\text{mon}$: find all docs containing any term ending with *mon*
 - Maintain an additional tree for terms *backwards*
 - Then retrieve all terms t in the range: $\text{nom} \leq t < \text{non}$
- Result: A set of terms that are matches for wildcard query
- Then retrieve documents that contain any of these terms

How to handle * in the middle of a term

- Example: c*sar
- We could look up c* and *sar in the B-tree and intersect the two term sets.
- Expensive
- Alternative: [permuterm](#) index
- Basic idea: Rotate every wildcard query, so that the * occurs at the end.
- Store each of these rotations in the dictionary, say, in a B-tree

Permuterm \rightarrow term mapping



Permuterm index

- For HELLO, we've stored: *hello\$, ello\$h, llo\$he, lo\$hel, and o\$hell*
- Queries
 - For X, look up X\$ (hello -> hello\$)
 - For X*, look up X*\$ (hel* -> hel*\$)
 - For *X, look up X\$* (*lo -> lo\$*)
 - For *X*, look up X* (*ll* -> ll*)
 - For X*Y, look up Y\$X* (hel*o -> o\$hel*)

Processing a lookup in the permuterm index

- Rotate query wildcard to the right
- Use B-tree lookup as before
- Problem: Permuterm more than **quadruples** the size of the dictionary compared to a regular B-tree. (empirical number)

k -gram indexes

- More space-efficient than permuterm index
- Enumerate all character k -grams (sequence of k characters) occurring in a term
- 2-grams are called **bigrams**.
- Example: from '*April is the cruelest month*' we get the bigrams: \$a ap pr ri il l\$ \$i is s\$ \$t th he e\$ \$c cr ru ue el le es st t\$ \$m mo on nt h\$
- \$ is a special word boundary symbol, as before.
- Maintain an inverted index from bigrams to the terms that contain the bigram

Postings list in a 3-gram inverted index



k -gram (bigram, trigram, . . .) indexes

- Note that we now have two different types of inverted indexes
- The term-document inverted index for finding documents based on a query consisting of terms
- The k -gram index for finding terms based on a query consisting of k -grams

Processing wildcarded terms in a bigram index

- Query mon^* can now be run as: \$m AND mo AND on
- Gets us all terms with the prefix *mon* . . .
- . . . but also many “false positives” like MOON.
- We must postfilter these terms against query.
- Surviving terms are then looked up in the term-document inverted index.
- k -gram index vs. permuterm index
 - k -gram index is more space efficient.
 - Permuterm index doesn’t require postfiltering.

gen* universit* - Google'da Ara - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://www.google.com.tr/#hl=tr&biw=1264&bih=815&q=gen*+universit*&aq=f&aqi=&aql=&oq=&fp=14706a139d059f65

Most Visited Getting Started Latest Headlines

gen* universit* - Google'da Ara

Web Görseller Haberler Çeviri Bloglar Gerçek zamanlı Gmail Diğer ▼

Google

gen* universit*

Yaklaşık 12.600.000 sonuç bulundu (0,21 saniye)

Gelişmiş arama

Her şey

Görseller

Videolar

Haberler

Daha fazla

Istanbul

Konumu değiştir

Web

Türkçe yazılmış sayfalar

Sayfaların bulunduğu ülke: Türkiye

Çevrilmiş sayfalar

Daha fazla arama aracı

Bunu mu demek istediniz? [gen* university*](#)

[Üniversiteler.gen.al: Üniversiteler, Devlet Üniversiteleri, Özel ...](#)

Üniversiteler.gen.al, Türkiye'deki devlet üniversiteleri ve özel üniversiteler hakkında bilgi edinebilmeniz için yayın yapıyor.

[www.universiteler.gen.al/](#) - Önbellek - Benzer

[MOLEKÜLER BİYOLOJİ ve GENETİK BÖLÜMÜ](#)

24 Şub 2011 ... İstanbul **Üniversitesi Fen** Fakültesi ... Bölümümüzün ERASMUS programı çerçevesinde Avrupadaki çeşitli **Üniversiteler ile** ikili anlaşmaları ...

[www.istanbul.edu.tr/fen/mol/](#) - Önbellek - Benzer

[Hacettepe Üniversitesi Genetik Ünitesi](#)

Hacettepe **Üniversitesi**. Tıp Fakültesi Çocuk Sağlığı ve Hastalıkları Ana Bilim Dalı Genetik Ünitesi & Çocuk Sağlığı Enstitüsü Temel Bilimler Ana Bilim Dalı ...

[www.gen.hacettepe.edu.tr/](#) - Önbellek - Benzer

[Boğaziçi Üniversitesi Moleküler Biyoloji ve Genetik Bölümü](#)

20 Oca 2011 ... Boğaziçi **Üniversitesi Moleküler** Biyoloji ve Genetik (MBG) bölümünde 7. çerçeve projesi kapsamında kurulan Biyogörüntüleme ve Analitik ...

[Dersler](#) - [İletişim](#) - [Mezunlar](#) - [Öğrenciler İçin Bilgiler](#)

[www.bio.boun.edu.tr/](#) - Önbellek - Benzer

Intention: you are looking for the University of Geneva, but don't know which accents to use for the French words for university and Geneva.

Google has very limited support for wildcard queries.

istanbul * university - Google'da Ara - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://www.google.com.tr/#hl=tr&source=hp&biw=1264&bih=815&q=istanbul+*+university&aq=f&aql=&oq=&fp=14706a139c

Most Visited Getting Started Latest Headlines

istanbul * university - Google'da Ara

Web Görseller Haberler Çeviri Bloglar Gerçek zamanlı Gmail Diğer ▼

Google

istanbul * university

Yaklaşık 936.000.000 sonuç bulundu (0,25 saniye)

Gelişmiş arama

Her şey

Görseller

Videolar

Haberler

Daha fazla

Istanbul

Konumu değiştir

Web

Türkçe yazılmış sayfalar

Sayfaların bulunduğu ülke: Türkiye

Çevrilmiş sayfalar

Daha fazla arama aracı

İstanbul Teknik Üniversitesi

İstanbul Teknik Üniversitesi, Akademik, İdari, Fakülteler ve Enstitüler, Akademik Takvim, Kampüsler ve Ulaşım, Araştırma.
www.itu.edu.tr/ - Önbellek - Benzer

İTÜ GIDA MÜHENDİSLİĞİ BÖLÜMÜ Ana Sayfa

İstanbul Teknik Üniversitesi Kimya Metalurji Fakültesi Gıda Mühendisliği ...
www.food.itu.edu.tr/ - Önbellek - Benzer

itu.edu.tr sitesinden daha fazla sonuç göster

İSTANBUL BİLGİ ÜNİVERSİTESİ

24 Şub 2011 ... İstanbul Bilgi University Official WEB-Site, İstanbul Bilgi Üniversitesi Resmi WEB Sitesi.
www.bilgi.edu.tr/ - Önbellek - Benzer

W OTEL İSTANBUL İşyeri Hekimi

İTÜ SÜREKLİ EĞİTİM MERKEZİ UZMANLIK SERTİFİKA PROGRAMLARI; İstanbul Institute Galatasaray Üniversitesi - Sağlık İletişimi Sertifika Programı ...
www.yenibiris.com/w_otel_istanbul/isyeri_hekimi/260376.ilan - Önbellek

İstanbul University - Wikipedia, the free encyclopedia - [Bu sayfanın çevirisini yap]
Black Sea Technical University (KTÜ) • İstanbul Technical University (İTÜ) • Middle East

According to Google search basics, 2010-04-29: "Note that the * operator works only on whole words, not parts of words."

m*nchen - Google'da Ara - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://www.google.com.tr/#hl=tr&biw=1264&bih=815&q=m*nchen&aq=f&aql=&oq=&fp=14706a139d059f65

Most Visited Getting Started Latest Headlines

m*nchen - Google'da Ara

Web Görseller Haberler Çeviri Bloglar Gerçek zamanlı Gmail Diğer ▼

Google

m*nchen

Yaklaşık 2.570.000 sonuç bulundu (0,17 saniye)

Gelişmiş arama

Her şey

- Görseller
- Videolar
- Haberler
- Daha fazla

İstanbul

Konumu değiştir

Web

Türkçe yazılmış sayfalar

Sayfaların bulunduğu ülke: Türkiye

Çevrilmiş sayfalar

Herhangi bir zaman

En yeni

Bunu mu demek istediniz? münchen En iyi 2 sonuç gösterildi

[muenchen.de - Munich Germany Bavaria bavarian german](#) - [Bu sayfanın çevirisini yap]

25 Feb 2011 ... 2011 Portal **München** Betriebs-GmbH & Co. KG Ein Service der Landeshauptstadt **München** und der Stadtwerke **München** GmbH. [Tourism - Traffic and Transport - City Life - Christmas Market](#) [www.muenchen.de/home/60093/Homepage.html](#) - Önbellek - Benzer

[Munich - Wikipedia, the free encyclopedia](#) - [Bu sayfanın çevirisini yap]

Munich (German: **München**, pronounced [ˈmʏnçən] (listen); Austro-Bavarian: Minga) is the capital city of Bavaria (Bayern), Germany. It is located on the ... [History - Geography - Demographics - Politics](#) [en.wikipedia.org/wiki/Munich](#) - Önbellek - Benzer

m*nchen için sonuçlar

[muenchen.de - Offizielles Stadtportal für München](#) - [Bu sayfanın çevirisini yap]

Jetzt zu M-Ökostrom wechseln und bis zu 60 € Wechselbonus kassieren. Hier online sparen! [www.swm.de/oekostrom](#) · Top Videos ... [Veranstaltungen & Tickets - Tourismus - Stadtplan - Stellenangebote, Jobs](#) [www.muenchen.de/](#) - Önbellek - Benzer

But this is not entirely true. Try [pythag*] and [m*nchen]

Why doesn't Google fully support wildcard queries?

Processing wildcard queries in the term-document index

- Problem 1: we must potentially execute a large number of Boolean queries.
 - Most straightforward semantics: Conjunction of disjunctions
 - For [gen* universit*]: geneva university OR geneva université OR genève university OR genève université OR general universities OR ...
 - Very expensive
- Problem 2: Users hate to type. If you encourage “laziness” people will respond!
- If abbreviated queries like [pyth* theo*] for [pythagoras’ theorem] are allowed, users will use them a lot.
- This would significantly increase the cost of answering queries.

Search

Type your search terms, use ‘*’ if you need to.
E.g., Alex* will match Alexander.

Spelling correction



Spelling correction

- Two principal uses
 - Correcting documents being indexed
 - Correcting user queries
- Two different methods for spelling correction
- **Isolated word** spelling correction
 - Check each word on its own for misspelling
 - Will not catch typos resulting in correctly spelled words, e.g., *I flew form Heathrow to Narita.*
- **Context-sensitive** spelling correction
 - Look at surrounding words
 - Can correct *form/from* error above

Document correction

- ▶ Especially needed for OCR'ed documents
 - ▶ Correction algorithms are tuned for this: rn/m
 - ▶ Can use domain-specific knowledge
 - ▶ E.g., OCR can confuse O and D more often than it would confuse O and I (adjacent on the QWERTY keyboard, so more likely interchanged in typing).
- ▶ But also: web pages and even printed material has typos
- ▶ But often we don't change the documents but aim to fix the query-document mapping



Query mis-spellings

- ▶ Our principal focus here
 - ▶ E.g., the query ***Albert Einstein***
- ▶ We can either
 - ▶ Retrieve documents indexed by the correct spelling, OR
 - ▶ Return several suggested alternative queries with the correct spelling
 - ▶ *Did you mean ... ?*





albert einstein



Ara

Yaklaşık 19.800.000 sonuç bulundu (0,15 saniye)

Gelişmiş arama

Her şey

Görseller

Videolar

Haberler

Daha fazla

İstanbul

Konumu değiştir

Web

Türkçe yazılmış sayfalar

Sayfaların bulunduğu

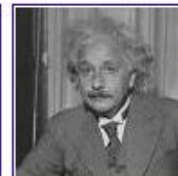
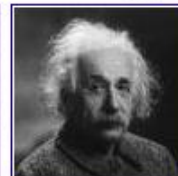
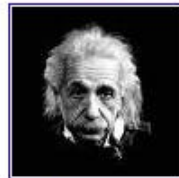
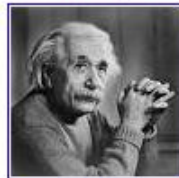
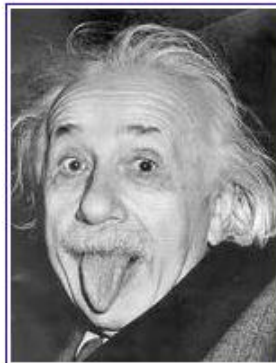
ülke: Türkiye

Çevrilmiş sayfalar

Tüm sonuçlar

Görselleri olan siteler

Daha fazla arama aracı

Bunu mu demek istediniz? [albert einstein](#)[albert einstein ile ilgili görseller](#) - Görseller hakkında kötüye kullanım bildirin[Albert Einstein - Vikipedi](#)

Albert Einstein (14 Mart 1879 - 18 Nisan 1955) , Yahudi asıllı Alman teorik ... 20. yüzyılın en önemli kuramsal fizikçisi olarak nitelenen **Albert Einstein**, ...
[tr.wikipedia.org/wiki/Albert_Einstein](#) - Önbellek - Benzer

[Albert Einstein - Wikipedia, the free encyclopedia](#) - [Bu sayfanın çevirisini yap]

Albert Einstein was a German-born theoretical physicist who discovered the ...

[en.wikipedia.org/wiki/Albert_Einstein](#) - Önbellek - Benzer

[wkipedia.org](#) sitesinden daha fazla sonuç göster

[Albert Einstein Hayatı \(Aralık Konuğu\) | EYLOS!...](#)

3 Ara 2007 ... **Albert Einstein** Hayatı (Aralık Konuğu) , Bilim , Windows değil EYLOS!...

[www.eylos.com](#) › Bilim - Önbellek - Benzer

Isolated word correction

- ▶ Fundamental premise – there is a lexicon from which the correct spellings come
- ▶ Two basic choices for this
 - ▶ A standard lexicon such as
 - ▶ Webster's English Dictionary
 - ▶ An “industry-specific” lexicon – hand-maintained
 - ▶ The lexicon of the indexed corpus
 - ▶ E.g., all words on the web
 - ▶ All names, acronyms etc.
 - ▶ (Including the mis-spellings)



Isolated word correction

- ▶ Given a lexicon and a character sequence Q , return the words in the lexicon closest to Q
- ▶ What's "closest"?
- ▶ We'll study several alternatives
 - ▶ Edit distance (Levenshtein distance)
 - ▶ Weighted edit distance
 - ▶ n -gram overlap

Edit distance

- The edit distance between string s_1 and string s_2 is the minimum number of basic operations that convert s_1 to s_2 .
- **Levenshtein distance:** The admissible basic operations are insert, delete, and replace (Edit distance usually refers to Levenshtein distance)
- Levenshtein distance *dog-do*: 1 (delete g)
- Levenshtein distance *cat-cart*: 1 (insert r)
- Levenshtein distance *cat-cut*: 1 (replace a with u)
- Levenshtein distance *cat-act*: 2 (replace c with a, replace a with c)
- **Damerau-Levenshtein distance** *cat-act*: 1 (transpose c with a)
- Damerau-Levenshtein includes transposition as a fourth possible operation.
- **Hamming distance:** only allows substitution (only applies to strings of the same length).



Levenshtein distance: Computation

		f	a	s	t
	0	1	2	3	4
c	1	1	2	3	4
a	2	2	1	2	3
t	3	3	2	2	2
s	4	4	3	2	3

Levenshtein distance: Algorithm

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), replace (cost 1), copy (cost 0)



Levenshtein distance: Algorithm

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), replace (cost 1), copy (cost 0)



Levenshtein distance: Algorithm

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), replace (cost 1), copy (cost 0)



Levenshtein distance: Algorithm

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), **replace (cost 1)**, copy (cost 0)



Levenshtein distance: Algorithm

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), replace (cost 1), **copy**
(cost 0)



Levenshtein distance: Example

		f	a	s	t
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
c	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>3</div><div>4</div><div>3</div><div>3</div></div>	<div><div>4</div><div>5</div><div>4</div><div>4</div></div>
a	<div><div>2</div><div>2</div></div>	<div><div>2</div><div>2</div><div>3</div><div>2</div></div>	<div><div>1</div><div>3</div><div>3</div><div>1</div></div>	<div><div>3</div><div>4</div><div>2</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>
t	<div><div>3</div><div>3</div></div>	<div><div>3</div><div>3</div><div>4</div><div>3</div></div>	<div><div>3</div><div>2</div><div>4</div><div>2</div></div>	<div><div>2</div><div>3</div><div>3</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>
s	<div><div>4</div><div>4</div></div>	<div><div>4</div><div>4</div><div>5</div><div>4</div></div>	<div><div>4</div><div>3</div><div>5</div><div>3</div></div>	<div><div>2</div><div>3</div><div>4</div><div>2</div></div>	<div><div>3</div><div>3</div><div>3</div><div>3</div></div>

Each cell of Levenshtein matrix

cost of getting here from my upper left neighbor (copy or replace)	cost of getting here from my upper neighbor (delete)
cost of getting here from my left neighbor (insert)	the minimum of the three possible “movements”; the cheapest way of getting here

Levenshtein distance: Example

		f	a	s	t
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
c	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>3</div><div>4</div><div>3</div><div>3</div></div>	<div><div>4</div><div>5</div><div>4</div><div>4</div></div>
a	<div><div>2</div><div>2</div></div>	<div><div>2</div><div>2</div><div>3</div><div>2</div></div>	<div><div>1</div><div>3</div><div>3</div><div>1</div></div>	<div><div>3</div><div>4</div><div>2</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>
t	<div><div>3</div><div>3</div></div>	<div><div>3</div><div>3</div><div>4</div><div>3</div></div>	<div><div>3</div><div>2</div><div>4</div><div>2</div></div>	<div><div>2</div><div>3</div><div>3</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>
s	<div><div>4</div><div>4</div></div>	<div><div>4</div><div>4</div><div>5</div><div>4</div></div>	<div><div>4</div><div>3</div><div>5</div><div>3</div></div>	<div><div>2</div><div>3</div><div>4</div><div>2</div></div>	<div><div>3</div><div>3</div><div>3</div><div>3</div></div>

Dynamic programming

- Optimal substructure: The optimal solution to the problem contains within it **subsolutions**, i.e., optimal solutions to subproblems.
- Overlapping subsolutions: The subsolutions overlap. These subsolutions are computed over and over again when computing the global optimal solution in a brute-force algorithm.
- Subproblem in the case of edit distance: what is the edit distance of two prefixes
- Overlapping subsolutions: We need most distances of prefixes 3 times – this corresponds to moving right, diagonally, down.

Weighted edit distance

- As above, but weight of an operation depends on the characters involved.
- Meant to capture keyboard errors, e.g., m more likely to be mistyped as n than as q .
- Therefore, replacing m by n is a smaller edit distance than by q .
- We now require a weight matrix as input.
- Modify dynamic programming to handle weights

Using edit distance for spelling correction

- Given query, first enumerate all character sequences within a preset (possibly weighted) edit distance
- Intersect this set with our list of “correct” words
- Then suggest terms in the intersection to the user.
- → exercise in a few slides

Exercise

- ① Compute Levenshtein distance matrix for OSLO – SNOW
- ② What are the Levenshtein editing operations that transform *cat* into *catcat*?

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
o	<div><div>1</div><div>1</div></div>				
s	<div><div>2</div><div>2</div></div>				
l	<div><div>3</div><div>3</div></div>				
o	<div><div>4</div><div>4</div></div>				

		s	n	o	w
	$\frac{\quad}{\quad} 0$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2} \frac{2}{?}$			
s	$\frac{2}{2}$				
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$			
s	$\frac{2}{2}$				
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				



		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{?}$		
s	$\frac{2}{2}$				
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$		
s	$\frac{2}{2}$				
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{\quad}{1} \frac{\quad}{1}$	$\frac{\quad}{2} \frac{\quad}{2}$	$\frac{\quad}{3} \frac{\quad}{3}$	$\frac{\quad}{4} \frac{\quad}{4}$
o	$\frac{\quad}{1} \frac{\quad}{1}$	$\frac{\quad}{1} \frac{\quad}{2} \frac{\quad}{2} \frac{\quad}{1}$	$\frac{\quad}{2} \frac{\quad}{3} \frac{\quad}{2} \frac{\quad}{2}$	$\frac{\quad}{2} \frac{\quad}{4} \frac{\quad}{3} \frac{\quad}{?}$	
s	$\frac{\quad}{2} \frac{\quad}{2}$				
l	$\frac{\quad}{3} \frac{\quad}{3}$				
o	$\frac{\quad}{4} \frac{\quad}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$	$\frac{2}{2}$	$\frac{2}{3}$	
s	$\frac{2}{2}$				
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{?}$
s	$\frac{2}{2}$				
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				



		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$				
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{?}$			
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$			
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{?}$		
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$		
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{?}$	
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{?}$
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$				
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{?}$			
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$			
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$	$\frac{2}{2}$	$\frac{2}{3}$	$\frac{4}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{3}{4}$
l	$\frac{3}{3}$	$\frac{3}{4}$	$\frac{2}{3}$		
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$		
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{4}$ $\frac{3}{?}$	
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{4}{3}$	
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{4}{3}$	$\frac{4}{4}$ $\frac{4}{?}$
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$				

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5}$ $\frac{3}{?}$			

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{4}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5}$ $\frac{3}{3}$			

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{3}{?}$		

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{3}{3}$		

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{2}$ $\frac{3}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{2}$ $\frac{4}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{3}$ $\frac{5}{3}$	$\frac{3}{3}$ $\frac{4}{3}$	$\frac{2}{4}$ $\frac{4}{?}$	

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{2}$ $\frac{3}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{2}$ $\frac{4}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{3}$ $\frac{5}{3}$	$\frac{3}{3}$ $\frac{4}{3}$	$\frac{2}{4}$ $\frac{4}{2}$	

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2} \frac{2}{1}$	$\frac{2}{2} \frac{3}{2}$	$\frac{2}{3} \frac{4}{2}$	$\frac{4}{3} \frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3} \frac{2}{1}$	$\frac{2}{2} \frac{3}{2}$	$\frac{3}{3} \frac{3}{3}$	$\frac{3}{4} \frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4} \frac{2}{2}$	$\frac{2}{3} \frac{3}{2}$	$\frac{3}{3} \frac{4}{3}$	$\frac{4}{4} \frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5} \frac{3}{3}$	$\frac{3}{4} \frac{3}{3}$	$\frac{2}{4} \frac{4}{2}$	$\frac{4}{3} \frac{5}{?}$

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{4}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{2}{4}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{2}{4}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{2}{4}$ $\frac{3}{2}$	$\frac{4}{5}$ $\frac{3}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{3}$ $\frac{2}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{2}{4}$ $\frac{4}{2}$	$\frac{4}{5}$ $\frac{3}{3}$

How do I read out the editing operations that transform OSLO into SNOW?



		s	n	o	w
	$\frac{\quad}{\quad} 0$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2} \frac{2}{1}$	$\frac{2}{2} \frac{3}{2}$	$\frac{2}{3} \frac{4}{2}$	$\frac{4}{3} \frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3} \frac{2}{1}$	$\frac{2}{2} \frac{3}{2}$	$\frac{3}{3} \frac{3}{3}$	$\frac{3}{4} \frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4} \frac{2}{2}$	$\frac{2}{3} \frac{3}{2}$	$\frac{3}{3} \frac{4}{3}$	$\frac{4}{4} \frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5} \frac{3}{3}$	$\frac{3}{4} \frac{3}{3}$	$\frac{2}{4} \frac{4}{2}$	$\frac{4}{3} \frac{5}{3}$

cost	operation	input	output
1	insert	*	w



		s	n	o	w
	$\begin{array}{ c } \hline 0 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 4 \\ \hline \end{array}$
o	$\begin{array}{ c } \hline 1 \\ \hline 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 2 \\ \hline 2 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 4 \\ \hline 3 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 5 \\ \hline 3 & 3 \\ \hline \end{array}$
s	$\begin{array}{ c } \hline 2 \\ \hline 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 2 \\ \hline 3 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 4 \\ \hline 4 & 3 \\ \hline \end{array}$
l	$\begin{array}{ c } \hline 3 \\ \hline 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 2 \\ \hline 4 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 3 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 4 \\ \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 4 \\ \hline 4 & 4 \\ \hline \end{array}$
o	$\begin{array}{ c } \hline 4 \\ \hline 4 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 3 \\ \hline 5 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline 4 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 4 \\ \hline 4 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 5 \\ \hline 3 & 3 \\ \hline \end{array}$

cost	operation	input	output
0	(copy)	o	o
1	insert	*	w

		s	n	o	w
	$\begin{array}{ c } \hline 0 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 4 \\ \hline \end{array}$
o	$\begin{array}{ c } \hline 1 \\ \hline 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 2 \\ \hline 2 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 4 \\ \hline 3 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 5 \\ \hline 3 & 3 \\ \hline \end{array}$
s	$\begin{array}{ c } \hline 2 \\ \hline 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 2 \\ \hline 3 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 4 \\ \hline 4 & 3 \\ \hline \end{array}$
l	$\begin{array}{ c } \hline 3 \\ \hline 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 2 \\ \hline 4 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 3 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 4 \\ \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 4 \\ \hline 4 & 4 \\ \hline \end{array}$
o	$\begin{array}{ c } \hline 4 \\ \hline 4 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 3 \\ \hline 5 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline 4 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 4 \\ \hline 4 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 5 \\ \hline 3 & 3 \\ \hline \end{array}$

cost	operation	input	output
1	replace	l	n
0	(copy)	o	o
1	insert	*	w



		s	n	o	w
	0	1 1	2 2	3 3	4 4
o	1 1	1 2 2 1	2 3 2 2	2 4 3 2	4 5 3 3
s	2 2	1 2 3 1	2 3 2 2	3 3 3 3	3 4 4 3
l	3 3	3 2 4 2	2 3 3 2	3 4 3 3	4 4 4 4
o	4 4	4 3 5 3	3 3 4 3	2 4 4 2	4 5 3 3

cost	operation	input	output
0	(copy)	s	s
1	replace	l	n
0	(copy)	o	o
1	insert	*	w

		s	n	o	w
	$\begin{array}{ c } \hline 0 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 4 \\ \hline \end{array}$
o	$\begin{array}{ c } \hline 1 \\ \hline 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 2 \\ \hline 2 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 4 \\ \hline 3 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 5 \\ \hline 3 & 3 \\ \hline \end{array}$
s	$\begin{array}{ c } \hline 2 \\ \hline 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 2 \\ \hline 3 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 4 \\ \hline 4 & 3 \\ \hline \end{array}$
l	$\begin{array}{ c } \hline 3 \\ \hline 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 2 \\ \hline 4 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 3 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 4 \\ \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 4 \\ \hline 4 & 4 \\ \hline \end{array}$
o	$\begin{array}{ c } \hline 4 \\ \hline 4 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 3 \\ \hline 5 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline 4 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 4 \\ \hline 4 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 5 \\ \hline 3 & 3 \\ \hline \end{array}$

cost	operation	input	output
1	delete	o	*
0	(copy)	s	s
1	replace	l	n
0	(copy)	o	o
1	insert	*	w

		c	a	t	c	a	t
	$\frac{\quad}{\quad} 0$	$\frac{1}{1} 1$	$\frac{2}{2} 2$	$\frac{3}{3} 3$	$\frac{4}{4} 4$	$\frac{5}{5} 5$	$\frac{6}{6} 6$
c	$\frac{\quad}{1} 1$	$\frac{0}{2} 2$	$\frac{2}{1} 3$	$\frac{3}{2} 4$	$\frac{3}{3} 5$	$\frac{5}{4} 6$	$\frac{6}{5} 7$
a	$\frac{\quad}{2} 2$	$\frac{2}{3} 1$	$\frac{0}{2} 2$	$\frac{2}{1} 3$	$\frac{3}{2} 4$	$\frac{3}{3} 5$	$\frac{5}{4} 6$
t	$\frac{\quad}{3} 3$	$\frac{3}{4} 2$	$\frac{2}{3} 1$	$\frac{0}{2} 2$	$\frac{2}{1} 3$	$\frac{3}{2} 4$	$\frac{3}{3} 5$

		c	a	t	c	a	t
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$	$\frac{5}{5}$	$\frac{6}{6}$
c	$\frac{1}{1}$	$\frac{0}{2}$	$\frac{2}{1}$	$\frac{3}{2}$	$\frac{3}{3}$	$\frac{5}{4}$	$\frac{6}{5}$
a	$\frac{2}{2}$	$\frac{2}{3}$	$\frac{0}{2}$	$\frac{2}{1}$	$\frac{3}{2}$	$\frac{3}{3}$	$\frac{5}{4}$
t	$\frac{3}{3}$	$\frac{3}{4}$	$\frac{2}{3}$	$\frac{0}{2}$	$\frac{2}{1}$	$\frac{3}{2}$	$\frac{3}{3}$

cost	operation	input	output
1	insert	*	c
1	insert	*	a
1	insert	*	t
0	(copy)	c	c
0	(copy)	a	a
0	(copy)	t	t

		c	a	t	c	a	t
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>	<div><div>5</div><div>5</div></div>	<div><div>6</div><div>6</div></div>
c	<div><div>1</div><div>1</div></div>	<div><div>0</div><div>2</div></div> <div><div>2</div><div>0</div></div>	<div><div>2</div><div>3</div></div> <div><div>1</div><div>1</div></div>	<div><div>3</div><div>4</div></div> <div><div>2</div><div>2</div></div>	<div><div>3</div><div>5</div></div> <div><div>3</div><div>3</div></div>	<div><div>5</div><div>6</div></div> <div><div>4</div><div>4</div></div>	<div><div>6</div><div>7</div></div> <div><div>5</div><div>5</div></div>
a	<div><div>2</div><div>2</div></div>	<div><div>2</div><div>1</div></div> <div><div>3</div><div>1</div></div>	<div><div>0</div><div>2</div></div> <div><div>2</div><div>0</div></div>	<div><div>2</div><div>3</div></div> <div><div>1</div><div>1</div></div>	<div><div>3</div><div>4</div></div> <div><div>2</div><div>2</div></div>	<div><div>3</div><div>5</div></div> <div><div>3</div><div>3</div></div>	<div><div>5</div><div>6</div></div> <div><div>4</div><div>4</div></div>
t	<div><div>3</div><div>3</div></div>	<div><div>3</div><div>2</div></div> <div><div>4</div><div>2</div></div>	<div><div>2</div><div>1</div></div> <div><div>3</div><div>1</div></div>	<div><div>0</div><div>2</div></div> <div><div>2</div><div>0</div></div>	<div><div>2</div><div>3</div></div> <div><div>1</div><div>1</div></div>	<div><div>3</div><div>4</div></div> <div><div>2</div><div>2</div></div>	<div><div>3</div><div>5</div></div> <div><div>3</div><div>3</div></div>

cost	operation	input	output
0	(copy)	c	c
1	insert	*	a
1	insert	*	t
1	insert	*	c
0	(copy)	a	a
0	(copy)	t	t

		c		a		t		c		a		t	
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{1}{1}$	$\frac{2}{1}$	$\frac{2}{1}$	$\frac{3}{2}$	$\frac{3}{2}$	$\frac{4}{3}$	$\frac{4}{3}$	$\frac{5}{4}$	$\frac{5}{4}$	$\frac{6}{5}$	$\frac{6}{5}$
c	$\frac{1}{1}$	$\frac{0}{2}$	$\frac{2}{0}$	$\frac{2}{1}$	$\frac{3}{1}$	$\frac{3}{2}$	$\frac{4}{2}$	$\frac{3}{3}$	$\frac{5}{3}$	$\frac{5}{4}$	$\frac{6}{4}$	$\frac{6}{5}$	$\frac{7}{5}$
a	$\frac{2}{2}$	$\frac{2}{3}$	$\frac{1}{1}$	$\frac{0}{2}$	$\frac{2}{1}$	$\frac{2}{1}$	$\frac{3}{1}$	$\frac{3}{2}$	$\frac{4}{2}$	$\frac{3}{3}$	$\frac{5}{3}$	$\frac{5}{4}$	$\frac{6}{4}$
t	$\frac{3}{3}$	$\frac{3}{4}$	$\frac{2}{2}$	$\frac{2}{3}$	$\frac{1}{1}$	$\frac{0}{2}$	$\frac{2}{0}$	$\frac{2}{1}$	$\frac{3}{1}$	$\frac{3}{2}$	$\frac{4}{2}$	$\frac{3}{3}$	$\frac{5}{3}$

cost	operation	input	output
0	(copy)	c	c
0	(copy)	a	a
1	insert	*	t
1	insert	*	c
1	insert	*	a
0	(copy)	t	t

		c		a		t		c		a		t	
		0	1 1	2 2	3 3	4 4	5 5	6 6					
c		1 1	0 2 2 0	2 3 1 1	3 4 2 2	3 5 3 3	5 6 4 4	6 7 5 5					
a		2 2	2 1 3 1	0 2 2 0	2 3 1 1	3 4 2 2	3 5 3 3	5 6 4 4					
t		3 3	3 2 4 2	2 1 3 1	0 2 2 0	2 3 1 1	3 4 2 2	3 5 3 3					

cost	operation	input	output
0	(copy)	c	c
0	(copy)	a	a
0	(copy)	t	t
1	insert	*	c
1	insert	*	a
1	insert	*	t



Spelling correction

- Now that we can compute edit distance: how to use it for isolated word spelling correction.
- k -gram indexes for isolated word spelling correction.
- Context-sensitive spelling correction
- General issues

Edit distance to all dictionary terms?

- ▶ Given a (mis-spelled) query – do we compute its edit distance to every dictionary term?
 - ▶ Expensive and slow
 - ▶ Alternative?
- ▶ How do we cut the set of candidate dictionary terms?
- ▶ One possibility is to use n -gram overlap for this
- ▶ This can also be used by itself for spelling correction.



n-gram overlap

- ▶ Enumerate all the *n*-grams in the query string as well as in the lexicon
- ▶ Use the *n*-gram index (recall wild-card search) to retrieve all lexicon terms matching any of the query *n*-grams
- ▶ Threshold by number of matching *n*-grams
 - ▶ Variants – weight by keyboard layout, etc.



Example with trigrams

- ▶ Suppose the text is **november**
 - ▶ Trigrams are *nov, ove, vem, emb, mbe, ber.*
- ▶ The query is **december**
 - ▶ Trigrams are *dec, ece, cem, emb, mbe, ber.*
- ▶ So 3 trigrams overlap (of 6 in each term)
- ▶ How can we turn this into a normalized measure of overlap?



One option – Jaccard coefficient

- ▶ A commonly-used measure of overlap
- ▶ Let X and Y be two sets; then the J.C. is

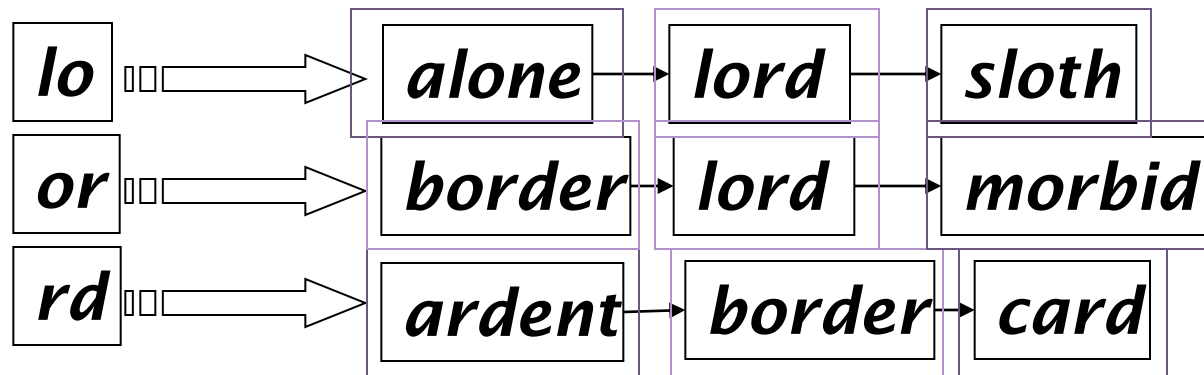
$$|X \cap Y| / |X \cup Y|$$

- ▶ Equals 1 when X and Y have the same elements and zero when they are disjoint
- ▶ X and Y don't have to be of the same size
- ▶ Always assigns a number between 0 and 1
 - ▶ Now threshold to decide if you have a match
 - ▶ E.g., if J.C. > 0.8, declare a match



Matching trigrams

- Consider the query **lord** – we wish to identify words matching 2 of its 3 bigrams (**lo**, **or**, **rd**)



Standard postings “merge” will enumerate ...

Context-sensitive spell correction

"flew form İstanbul Ataturk Airport" - Google'da Ara - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://www.google.com.tr/search?hl=tr&client=firefox-a&hs=pCu&rls=org.mozilla%3Aen-US%3Aofficial&channel=s&q="flew+form+İst

Most Visited Getting Started Latest Headlines

"flew form İstanbul Ataturk Airport" ...

Web Görseller Haberler Çeviri Bloglar Gerçek zamanlı Gmail Diğer ▼

Google

"flew form İstanbul Ataturk Airport"

Yaklaşık 0 sonuç bulundu (0,26 saniye)

Gelişmiş arama

Her şey

Görseller

Videolar

Haberler

Daha fazla

İstanbul

Konumu değiştir

Web

Türkçe yazılmış sayfalar

Sayfaların bulunduğu ülke: Türkiye

Çevrilmiş sayfalar

Daha fazla arama aracı

Bunu mu demek istediniz? "flew **from** İstanbul Ataturk Airport" En iyi 2 sonuç gösterildi

port | Port Survey

We flew from İstanbul Atatürk Airport. Can you share your impression of international terminals as a place where people from all over the world come ...

www.tavnewsport.com/Lale-Mansur_491/port - Önbellek

Mount Nemrut...where to go? - Lonely Planet travel forum

5 gönderi - 4 yazar - Son gönderi: 12 Tem 2010

hi there.....last year I flew from İstanbul Ataturk airport with Turkish airlines to Malayta....it was a very early morning ...

www.lonelyplanet.com/thorntree/thread.jspa?threadID=1924825 - Önbellek

"flew form İstanbul Ataturk Airport" için sonuçlar

Aradığınız - "flew form İstanbul Ataturk Airport" - ile ilgili hiçbir arama sonucu mevcut değil.

Öneriler:

Tüm kelimeleri doğru yazdığınızdan emin olun.

Context-sensitive correction

- ▶ Need surrounding context to catch this.
- ▶ First idea: retrieve dictionary terms close (in weighted edit distance) to each query term
- ▶ Now try all possible resulting phrases with one word “fixed” at a time
 - ▶ *flew from Istanbul Ataturk Airport*
 - ▶ *fled form Istanbul Ataturk Airport*
 - ▶ *flea form Istanbul Ataturk Airport*
- ▶ **Hit-based spelling correction:** Suggest the alternative that has lots of hits.



Exercise

- ▶ Suppose that for “***flew form Istanbul Ataturk Airport***” we have 7 alternatives for flew, 20 for form, 3 for Istanbul, 2 for Ataturk, and 3 for airport.

How many “corrected” phrases will we enumerate in this scheme?

Another approach

- ▶ Break phrase query into a conjunction of biwords (Lecture 2).
- ▶ Look for biwords that need only one term corrected.
- ▶ Enumerate phrase matches and ... rank them!

General issues in spell correction

- ▶ We enumerate multiple alternatives for “Did you mean?”
- ▶ Need to figure out which to present to the user
- ▶ Use heuristics
 - ▶ The alternative hitting most docs
 - ▶ Query log analysis + tweaking
 - ▶ For especially popular, topical queries
- ▶ Spell-correction is computationally expensive
 - ▶ Avoid running routinely on every query?
 - ▶ Run only on queries that matched few docs





Soundex



Soundex

- ▶ Class of heuristics to expand a query into **phonetic** equivalents
 - ▶ Language specific – mainly for names
 - ▶ E.g., ***chebyshev*** → ***tchebycheff***

Soundex – typical algorithm

- ▶ Turn every token to be indexed into a 4-character reduced form
- ▶ Do the same with query terms
- ▶ Build and search an index on the reduced forms
 - ▶ (when the query calls for a soundex match)

Soundex – typical algorithm

1. Retain the first letter of the word.
2. Change all occurrences of the following letters to '0' (zero):
'A', 'E', 'I', 'O', 'U', 'H', 'W', 'Y'.
3. Change letters to digits as follows:
 - ▶ B, F, P, V → 1
 - ▶ C, G, J, K, Q, S, X, Z → 2
 - ▶ D, T → 3
 - ▶ L → 4
 - ▶ M, N → 5
 - ▶ R → 6
4. Remove all pairs of consecutive digits.
5. Remove all zeros from the resulting string.
6. Pad the resulting string with trailing zeros and return the first four positions, which will be of the form <uppercase letter> <digit> <digit> <digit>.



Example: Soundex of *HERMAN*

- Retain H
- *ERMAN* → *ORMON*
- *ORMON* → 06505
- 06505 → 06505
- 06505 → 655
- Return *H655*
- Note: *HERMANN* will generate the same code

Soundex

- ▶ Soundex is the classic algorithm, provided by most databases (Oracle, Microsoft, ...)
- ▶ How useful is soundex?
- ▶ Not very – for information retrieval
- ▶ Okay for “high recall” tasks, though biased to names of certain nationalities

What queries can we process?

- ▶ We have
 - ▶ Positional inverted index with skip pointers
 - ▶ Wild-card index
 - ▶ Spell-correction
 - ▶ Soundex
- ▶ Queries such as
(SPELL(moriset) /3 toron*to) OR SOUNDEX(chaikofski)



References

- ▶ *Introduction to Information Retrieval*, chapter 3
- ▶ The slides were adapted from the book's companion website:
 - ▶ <http://nlp.stanford.edu/IR-book/information-retrieval-book.html>
- ▶ A nice example and an applet for edit distance.
 - ▶ <http://www.merriampark.com/ld.htm>
- ▶ Nice reading on spell correction:
 - ▶ Peter Norvig: How to write a spelling corrector
<http://norvig.com/spell-correct.html>
- ▶ Soundex Algorithm demo:
 - ▶ <http://www.creativyst.com/Doc/Articles/SoundExI/SoundExI.htm#Top>

