

SEZNAM.CZ

**Náš život s
jazykovými modely:
,Sumarizace
dokumentu‘**

Martin Bachura



SEZNAM.CZ

Vektorová reprezentace dokumentu

Martin Bachura



Text / Document summarization

Text extraction vs Text generation

Vectors

Research vs Engineering

Model designs vs industrial applications

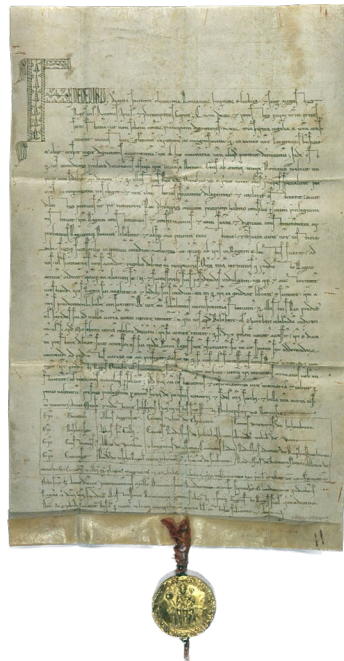
Evaluation

GPT vs Smaller LM models



Our problem

Retrieve appropriate documents /
websites from index for given
query



SEZNAM.CZ

jak se nazýval dokument z roku 1212, který mj. zaručoval českým panovníkům ×

Internet Obrázky Zboží Mapy Video Zprávy Firmy Slovník

Zlatá bula sicilská – Wikipedie
cs.wikipedia.org/wiki/zlat%C3%A1-bula-sicilsk%C3%A1
Název listiny, resp. listin, je odvozen od pečeti, která je k **dokumentu** přivěšena. Fridrich II. jako král Sicílie disponoval tehdy pouze pečeti tohoto království.
[Obsah prvního privilegia](#) [Poznámky](#) [Literatura](#) [Související články](#) [Externí odkazy](#)

Přemyslovci – Wikipedie
cs.wikipedia.org/wiki/přemyslovci
Ve vedlejší (levobočné) opavské linii Přemyslovci vymřeli (po meči) až **roku** 1521.
[Původ dynastie](#) [Název](#) [Období vlády](#) [Poslední Přemyslovci](#) [Tělesná charakteristika](#)

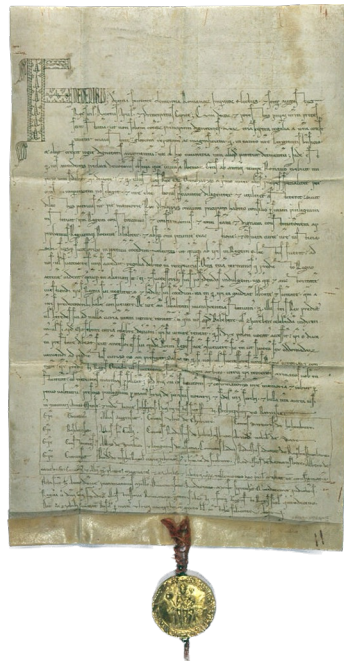
700. výročí vymření Přemyslovců po meči Václavem III
zlate-mince.cz/crs_2006_premyslovci_info.htm
700. výročí vymření Přemyslovců po meči Václavem III, stříbrná pamětní mince **České** národní banky v hodnotě 200 Kč. Získáte ji ve specializované prodejně numismatiky v Obecním domě v Praze.

Příběhy výzkumu a vývoje: Jak rychle najít Zlatou bulu...
tripartita.cz/pribehy-vyzkumu-a-vyvoje-jak-rychle-najit-zlatou...
Tripartita, oficiálním názvem Rada hospodářské a sociální dohody **České** republiky (RHSD ČR) je společným dobrovolným dohadovacím a iniciativním orgánem odborů, zaměstnavatelů a vlády



Our problem

Retrieve appropriate documents /
websites from **vector** index for
given query



SEZNAM.CZ

jak se nazýval dokument z roku 1212, který mj. zaručoval českým panovníkům ×

Internet Obrázky Zboží Mapy Video Zprávy Firmy Slovník

Zlatá bula sicilská – Wikipedie
cs.wikipedia.org/wiki/zlat%C3%A1-bula-sicilsk%C3%A1
Název listiny, resp. listin, je odvozen od pečeti, která je k **dokumentu** přivěšena. Fridrich II. jako král Sicílie disponoval tehdy pouze pečeti tohoto království.
[Obsah prvního privilegia](#) [Poznámky](#) [Literatura](#) [Související články](#) [Externí odkazy](#)

Přemyslovci – Wikipedie
cs.wikipedia.org/wiki/přemyslovci
Ve vedlejší (levobočné) opavské linii Přemyslovci vymřeli (po meči) až **roku** 1521.
[Původ dynastie](#) [Název](#) [Období vlády](#) [Poslední Přemyslovci](#) [Tělesná charakteristika](#)

700. výročí vymření Přemyslovců po meči Václavem III
zlate-mince.cz/crs_2006_premyslovci_info.htm
700. výročí vymření Přemyslovců po meči Václavem III, stříbrná pamětní mince **České** národní banky v hodnotě 200 Kč. Získáte ji ve specializované prodejně numismatiky v Obecním domě v Praze.

Příběhy výzkumu a vývoje: Jak rychle najít Zlatou bulu...
tripartita.cz/pribehy-vyzkumu-a-vyvoje-jak-rychle-najit-zlatou...
Tripartita, oficiálním názvem Rada hospodářské a sociální dohody **České** republiky (RHSD ČR) je společným dobrovolným dohadovacím a iniciativním orgánem odborů, zaměstnavatelů a vlády



Our problem

Jak se nazýval dokument z roku 1212 který mj. zajišťoval českým panovníkům dědičný královský titul

Zlatá bula sicilská (lat. *Bulla Aurea Siciliae*) je označení pro jednu hlavní nebo soubor tří souvisejících listin, které vydal 26. září 1212 v Basileji budoucí římský král Fridrich II. českému Přemyslu Otakarovi I. Listina především obsahuje privilegium, které potvrzuje dědičnost titulu českého krále a stanovuje jeho volbu na domácí půdě a to, že římský panovník bude volbu jen stvrzovat.

Přemysl Otakar ji získal jako odměnu za podporu, již Fridrichovi poskytl pro volbu římským králem. Novodobý název listiny, resp. listin, je odvozen od zlaté buly (kovové pečeti) sicilského krále, která je k dokumentu přivěšena. Fridrich v té době ještě neměl takto slavnostní pečeť římského krále, použil svou pečeť krále jižní Itálie.^[1]

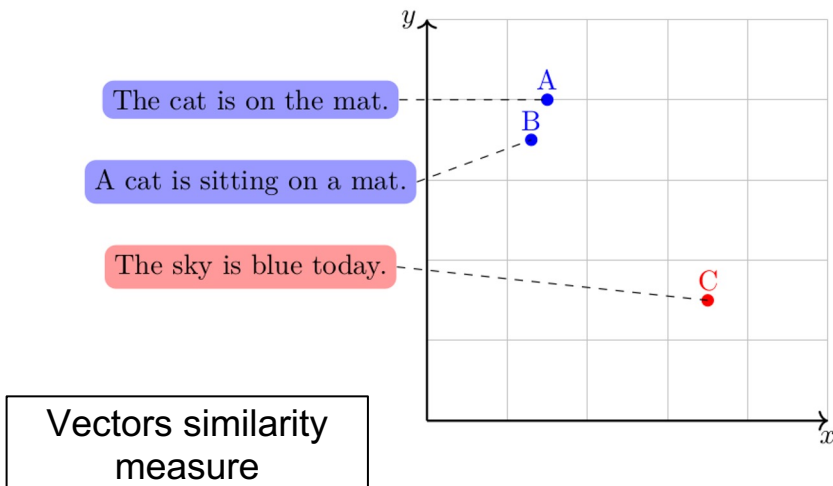
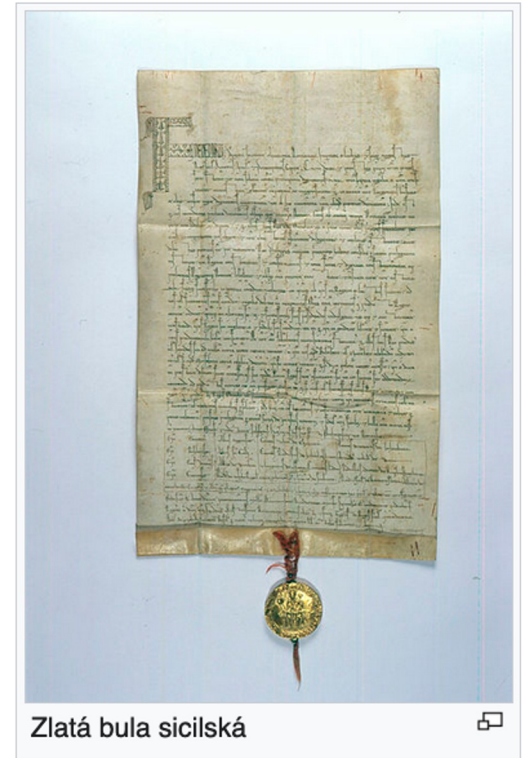


Our problem

Jak se nazýval dokument z roku 1212 který mj. zajišťoval českým panovníkům dědičný královský titul

Zlatá bula sicilská (lat. *Bulla Aurea Siciliae*) je označení pro jednu hlavní nebo soubor tří souvisejících listin, které vydal 26. září 1212 v Basileji budoucí římský král Fridrich II. českému Přemyslu Otakarovi I. Listina především obsahuje privilegium, které potvrzuje dědičnost titulu českého krále a stanovuje jeho volbu na domácí půdě a to, že římský panovník bude volbu jen stvrzovat.

Přemysl Otakar ji získal jako odměnu za podporu, již Fridrichovi poskytl pro volbu římským králem. Novodobý název listiny, resp. listin, je odvozen od zlaté buly (kovové pečeti) sicilského krále, která je k dokumentu přivěšena. Fridrich v té době ještě neměl takto slavnostní pečeť římského krále, použil svou pečeť krále jižní Itálie.^[1]



Our problem

Jak se nazýval dokument z roku 1212 který mj. zajišťoval českým panovníkům dědičný královský titul

Zlatá bula sicilská (lat. *Bulla Aurea Siciliae*) je označení pro jednu hlavní nebo soubor tří souvisejících **listin**, které vydal **26. září 1212** v **Basileji** budoucí římský král **Fridrich II.** českému **Přemyslu Otakarovi I.** Listina především obsahuje **privilegium**, které potvrzuje dědičnost titulu **českého krále** a stanovuje jeho volbu na domácí půdě a to, že **římský panovník** bude volbu jen stvrzovat.

Přemysl Otakar ji získal jako odměnu za podporu, již Fridrichovi poskytl pro volbu římským králem.

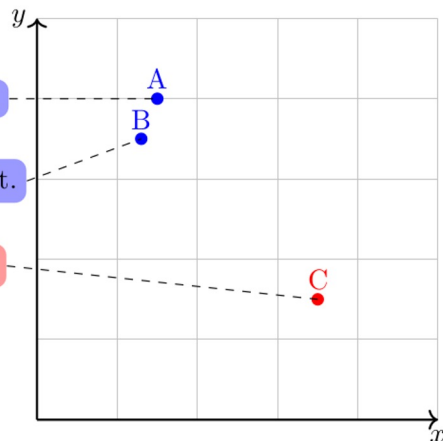
Novodobý název listiny, resp. listin, je odvozen od **zlaté buly** (kovové pečeti) **sicilského krále**, která je k dokumentu přivěšena. Fridrich v té době ještě neměl takto slavnostní pečeť římského krále, použil svou pečeť krále jižní Itálie.^[1]



The cat is on the mat.

A cat is sitting on a mat.

The sky is blue today.



Vectors similarity measure

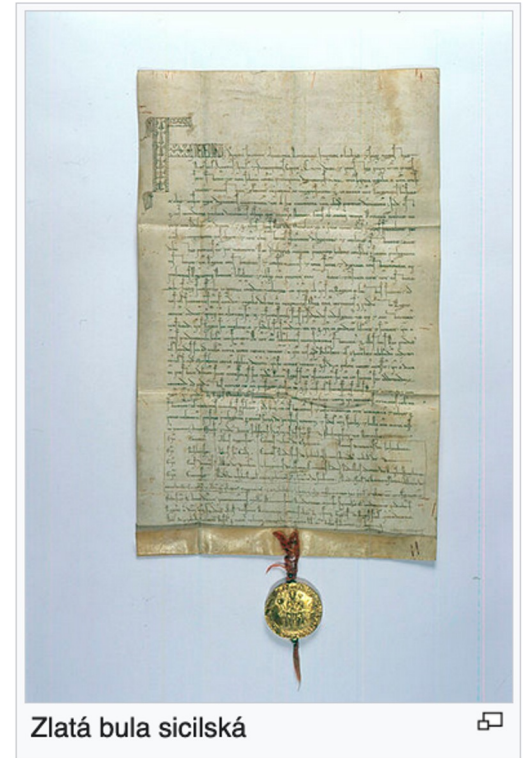


Our problem

Jak se nazýval dokument z roku 1212 který mj. zajišťoval českým panovníkům dědičný královský titul

Zlatá bula sicilská (lat. *Bulla Aurea Siciliae*) je označení pro jednu hlavní nebo soubor tří souvisejících **listin**, které vydal **26. září 1212** v **Basileji** budoucí římský král **Fridrich II.** českému **Přemyslu Otakarovi I.** Listina především obsahuje **privilegium**, které potvrzuje dědičnost titulu **českého krále** a stanovuje jeho volbu na domácí půdě a to, že **římský panovník** bude volbu jen stvrzovat.

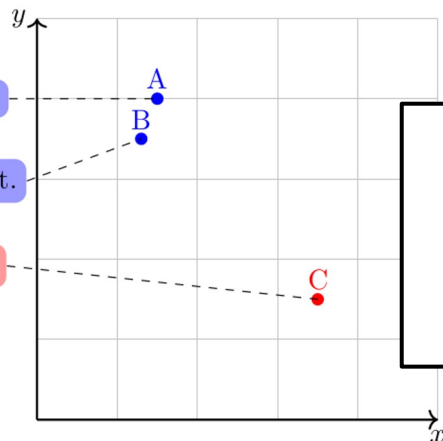
Přemysl Otakar ji získal jako odměnu za podporu, již Fridrichovi poskytl pro volbu římským králem. Novodobý název listiny, resp. listin, je odvozen od **zlaté buly** (kovové pečeti) **sicilského krále**, která je k dokumentu přivěšena. Fridrich v té době ještě neměl takto slavnostní pečeť římského krále, použil svou



The cat is on the mat.

A cat is sitting on a mat.

The sky is blue today.



Vectorization of proper text segments: selection of proper segments

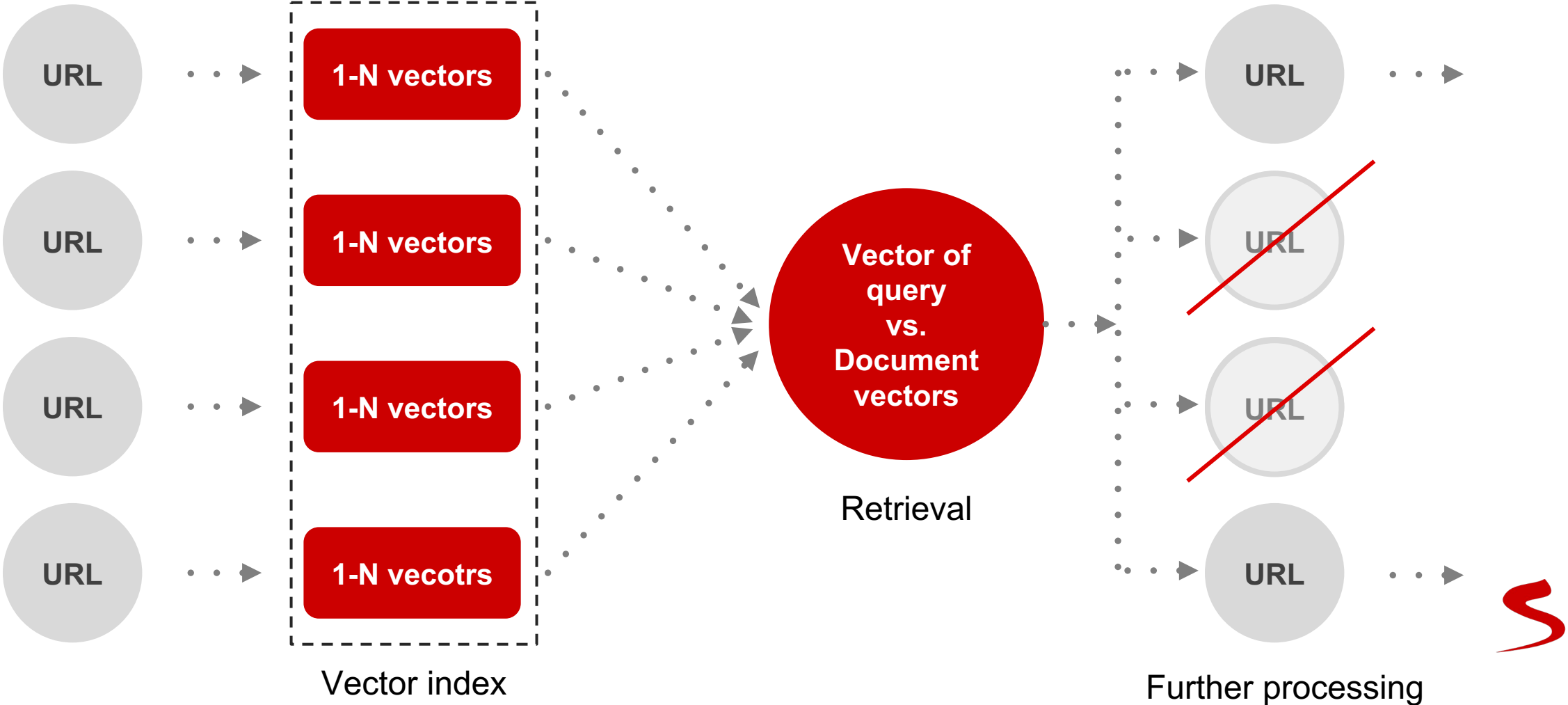
Vectors similarity measure



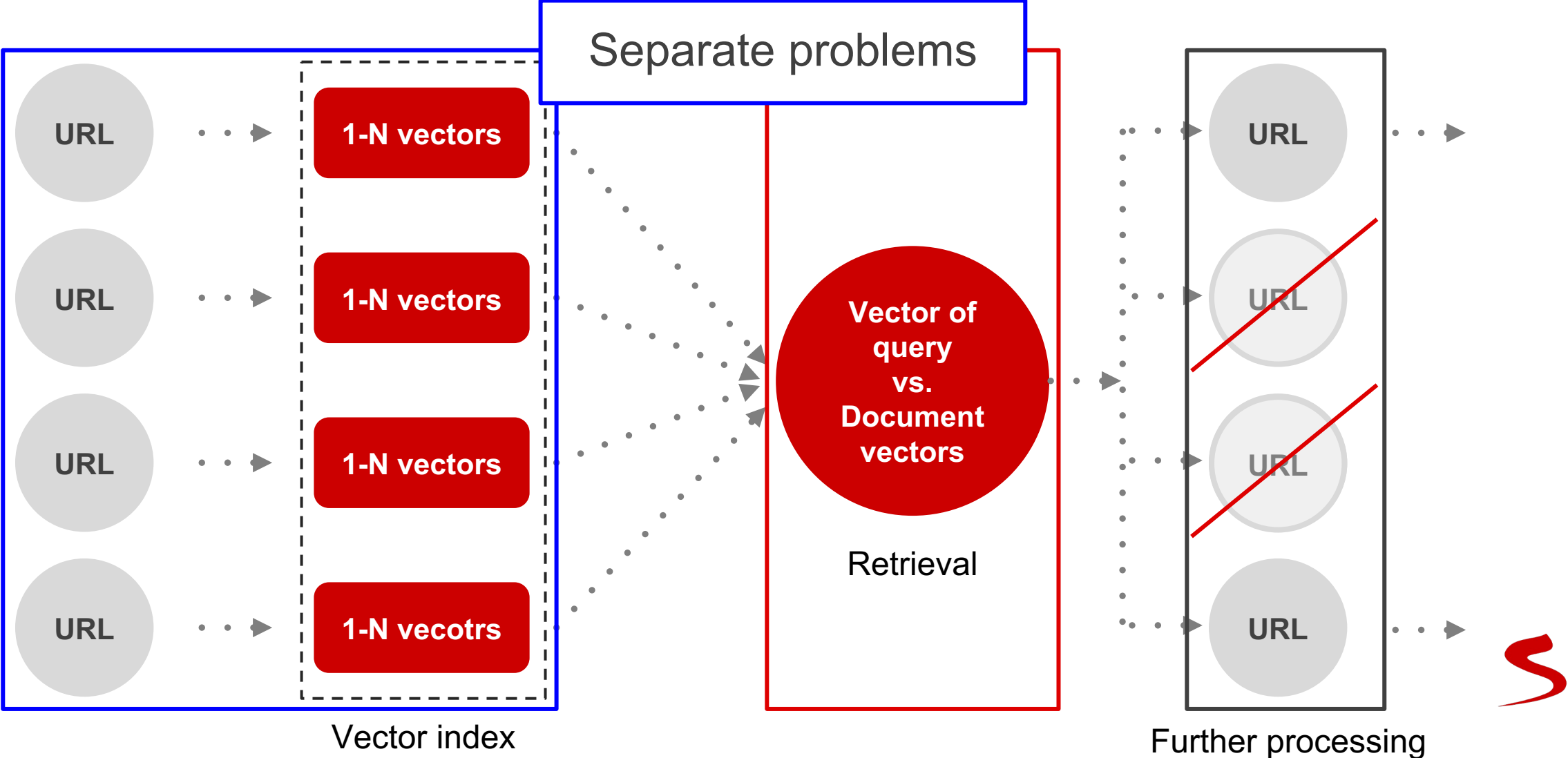
Vector branch of Seznam search



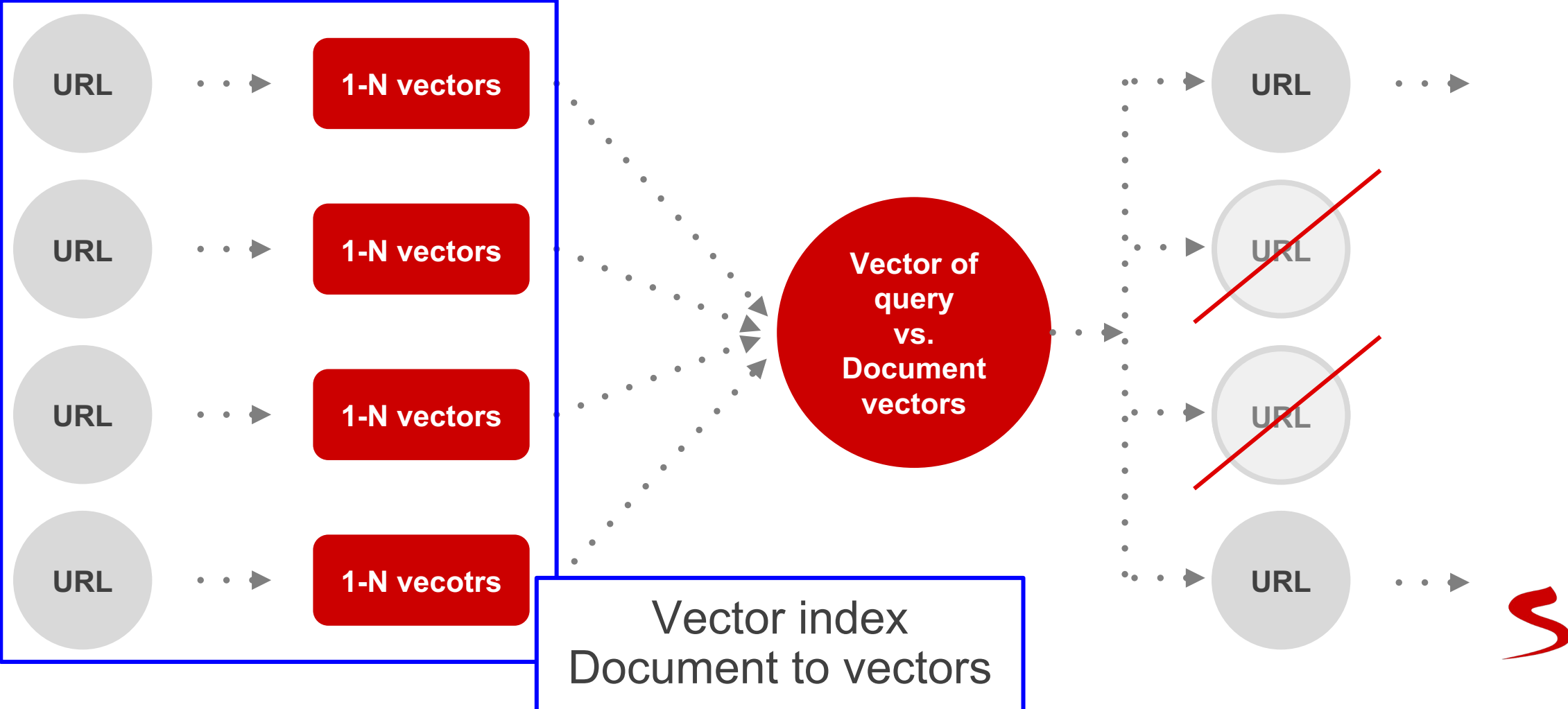
Vector branch of our search engine



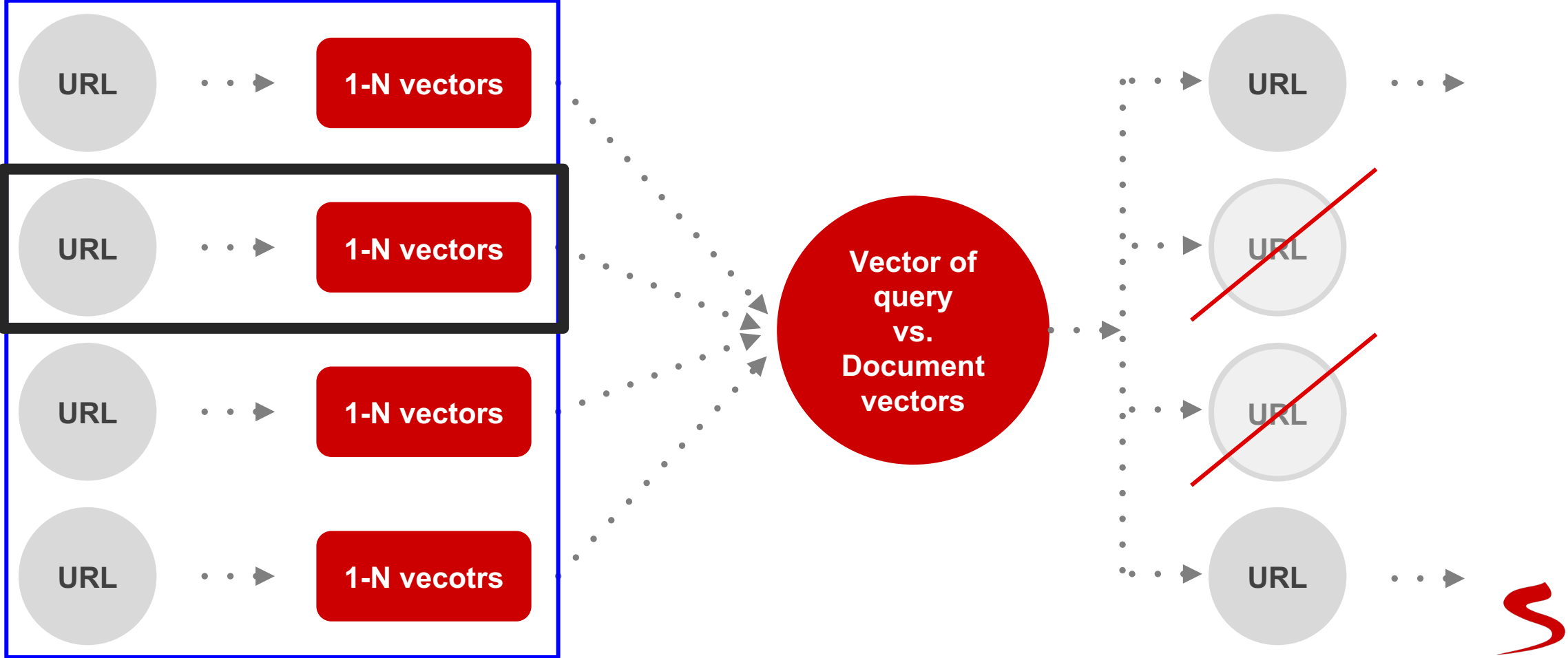
Vector branch of our search engine



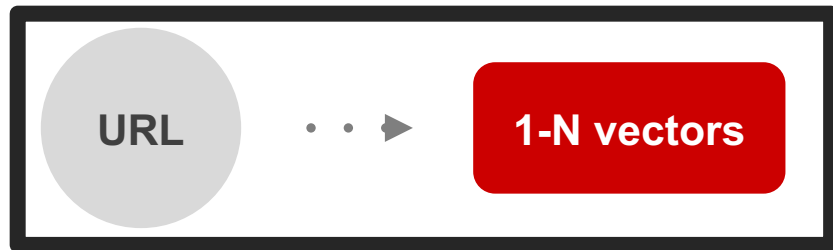
Vector branch of our search engine - our problem



Vector branch of our search engine - our problem



Our problem



Vectorization of proper text segments: selection of proper segments

Zlatá bula sicilská (lat. *Bulla Aurea Siciliae*) je označení pro jednu hlavní nebo soubor tří souvisejících **listin**, které vydal **26. září 1212** v **Basileji** budoucí římský král **Fridrich II.** českému **Přemyslu Otakarovi I.** Listina především obsahuje **privilegium**, které potvrzuje dědičnost titulu **českého krále** a stanovuje jeho volbu na domácí půdě a to, že **římský panovník** bude volbu jen stvrzovat.

Přemysl Otakar ji získal jako odměnu za podporu, již Fridrichovi poskytl pro volbu římským králem.

Novodobý název listiny, resp. listin, je odvozen od **zlaté buly** (kovové pečeti) **sicilského krále**, která je k dokumentu přivěšena. Fridrich v té době ještě neměl takto slavnostní pečeť římského krále, použil svou pečeť krále jižní Itálie.^[1]



Vectors



Vectors - properties

- Storage limitations - how many vectors?
- >> 1 document, whole Czech internet, refreshed
- Limited computational power and \$\$
- Documents are large
- Retrieval task - one vector one topic, reasonable vector space properties
- Vectors from text, texts - interpretability



Vectors - properties



- Storage limitations - how many vectors?
- >> 1 document, whole Czech internet, refreshed
- Limited computational power and \$\$
- Documents are large
- Retrieval task - one vector one topic, reasonable vector space properties
- Vectors from text, texts - interpretability



Vectors - properties



- Unlimited storage - up to N vectors per document
- >> 1 document, whole Czech internet, refreshed
- Limited computational power and \$\$
- Documents are large
- Retrieval task - one vector one topic, reasonable vector space properties
- Vectors from text, texts - interpretability



Vectors - properties



- Unlimited storage - up to N vectors per document
- >> 1 document, whole Czech internet, refreshed
- Limited computational power and \$\$
- Documents are large
- Retrieval task - one vector one topic, reasonable vector space properties
- Vectors from text, texts - interpretability



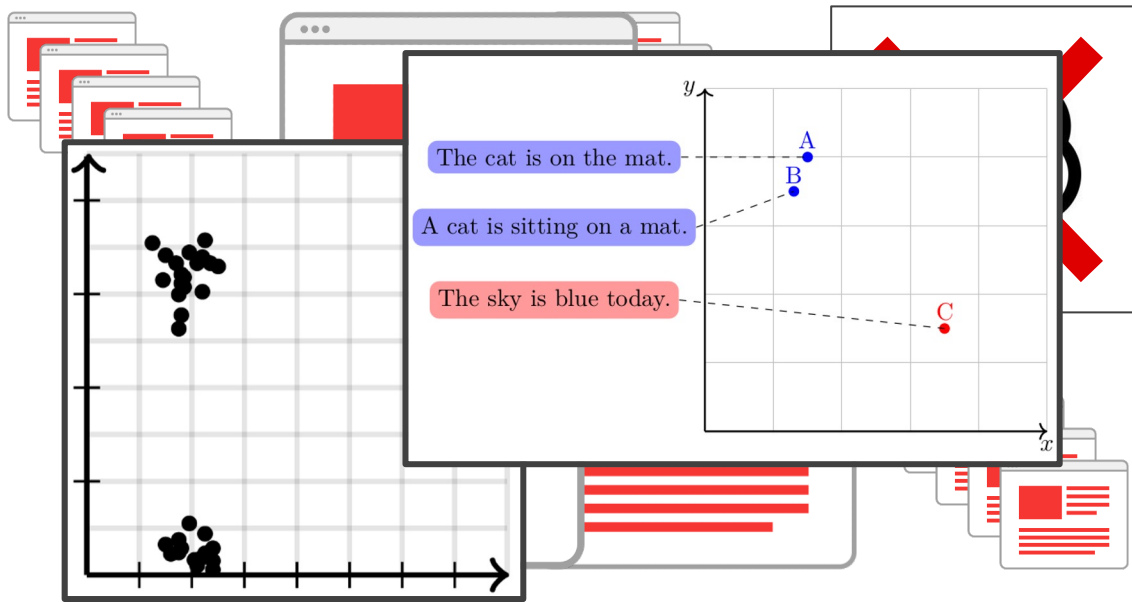
Vectors - properties



- Unlimited storage - up to N vectors per document
- >> 1 document, whole Czech internet, refreshed
- Limited computational power and \$\$
- Documents are large
- Retrieval task - one vector one topic, reasonable vector space properties
- Vectors from text, texts - interpretability



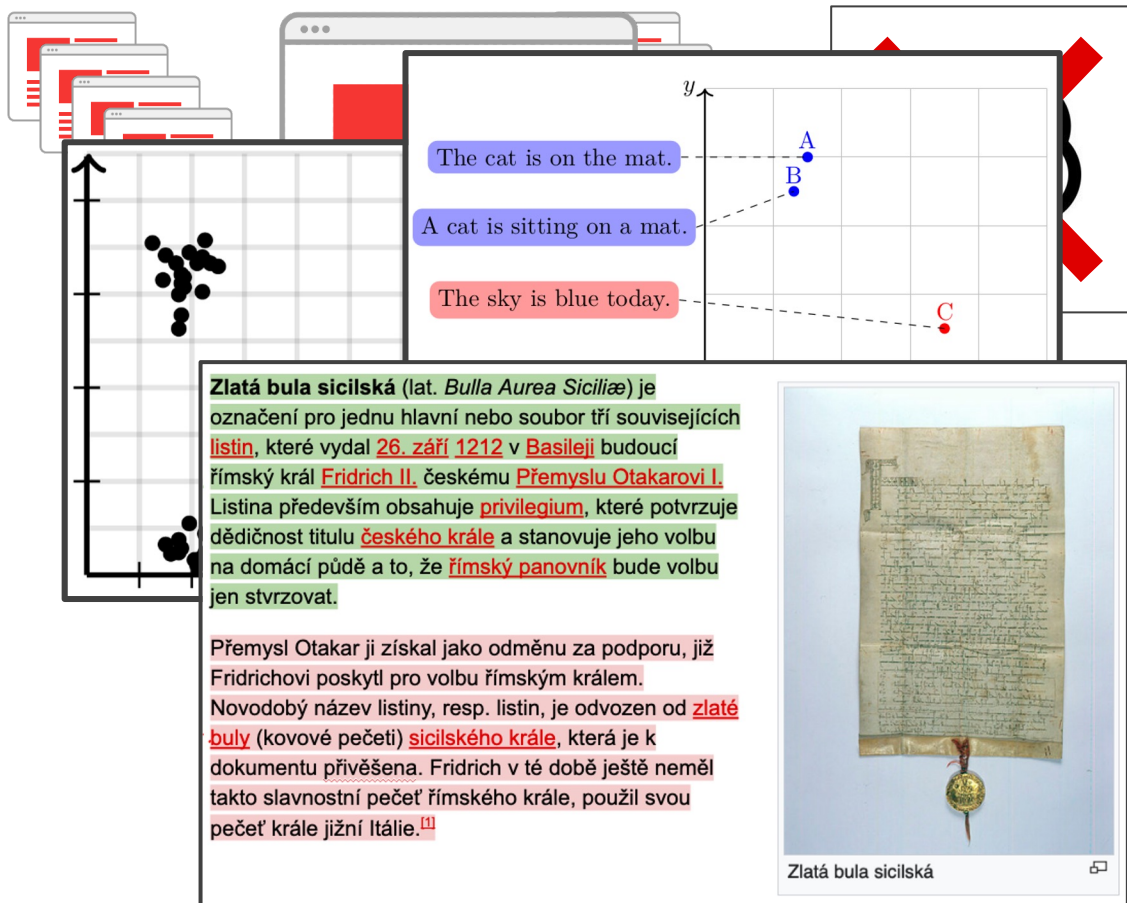
Vectors - properties



- Unlimited storage - up to N vectors per document
- >> 1 document, whole Czech internet, refreshed
- Limited computational power and \$\$
- Documents are large
- Retrieval task - one vector one topic, reasonable vector space properties
- Vectors from text, texts - interpretability



Vectors - properties



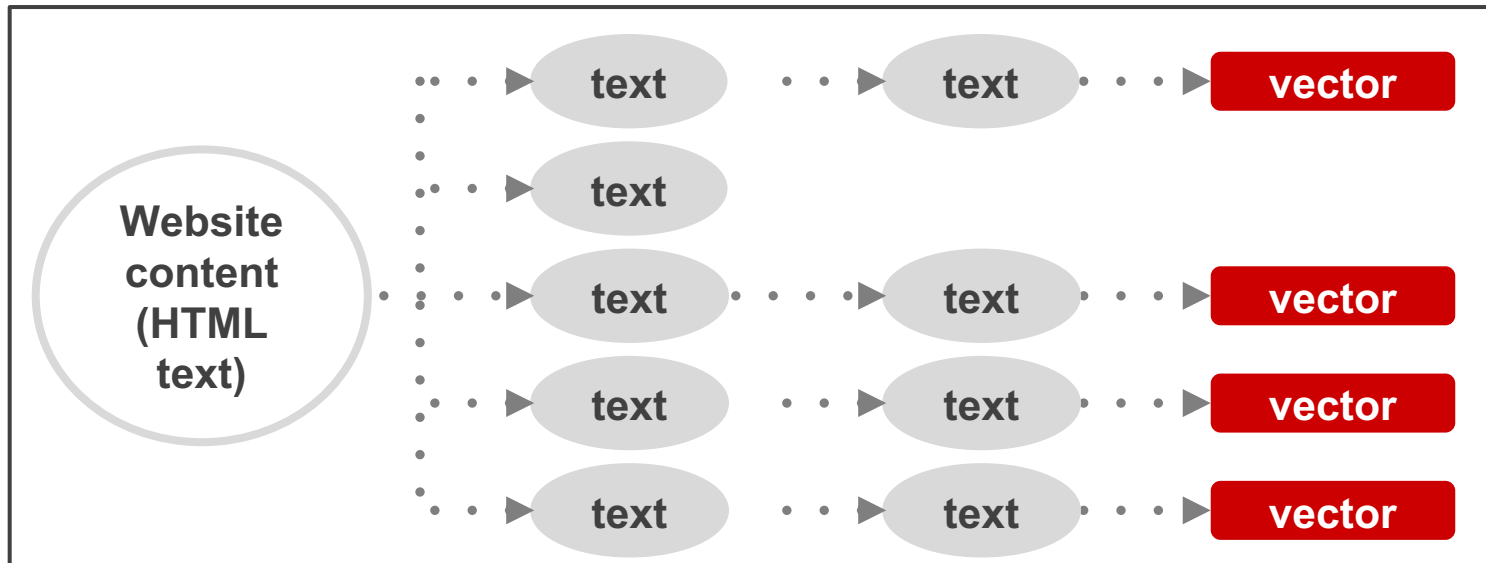
- Unlimited storage - up to N vectors per document
- >> 1 document, whole Czech internet, refreshed
- Limited computational power and \$\$
- Documents are large
- Retrieval task - one vector one topic, reasonable vector space properties
- Vectors from text, texts - interpretability



Our research



Our solution



Works in context of vector search

Applicable to production

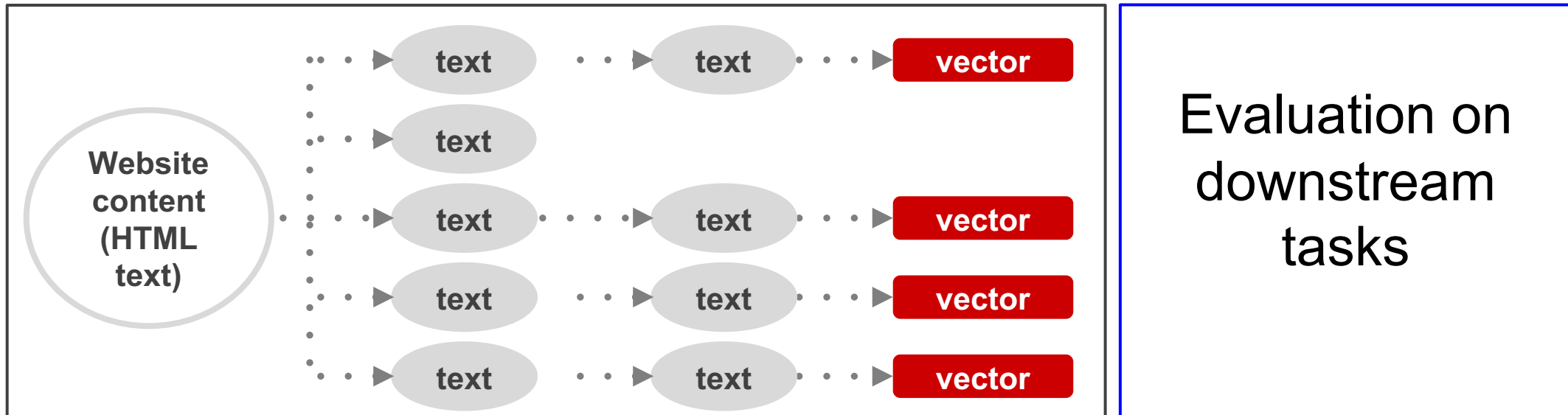
'Single topic' vectors

Rather smaller number of vectors

Split - Extract - Vectorize



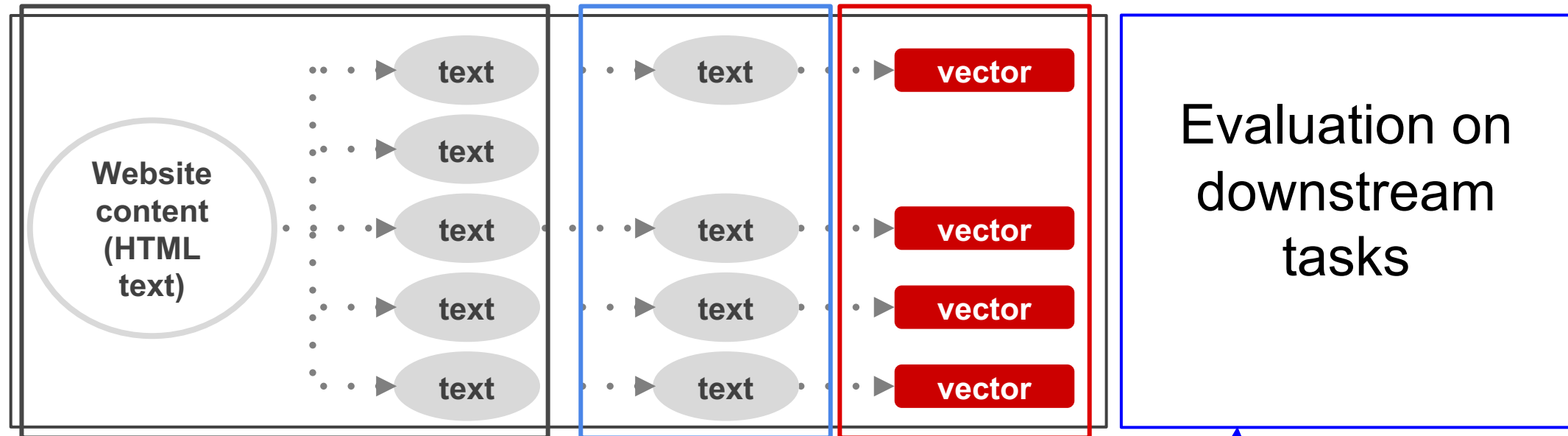
Our research



Split - Extract - Vectorize - Evaluate



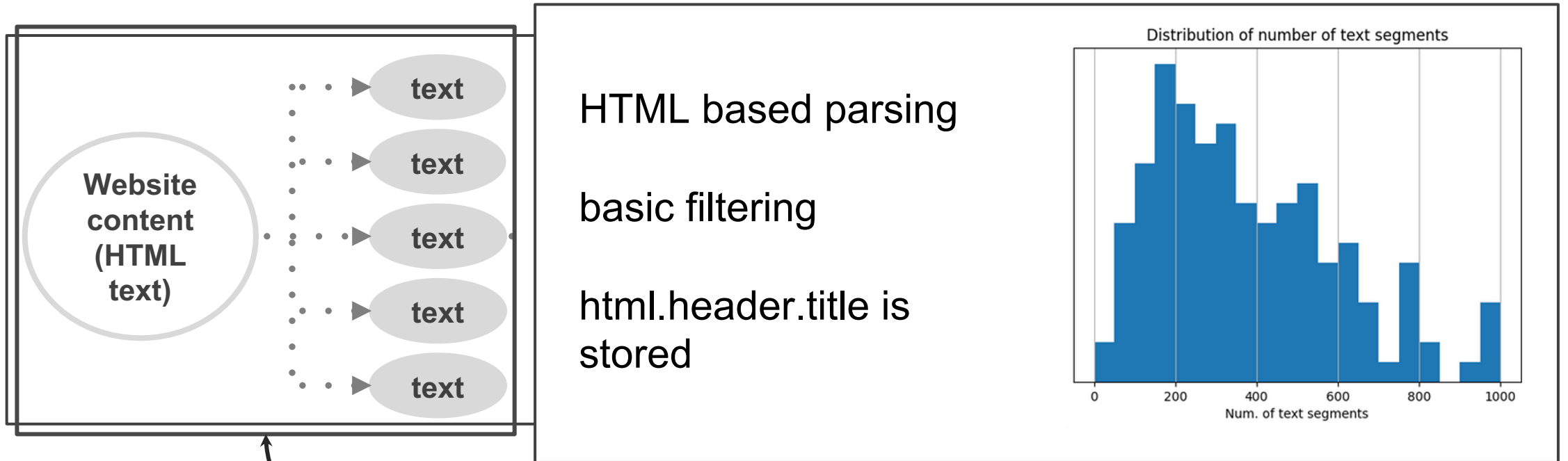
Our research



Split - Extract - Vectorize - Evaluate



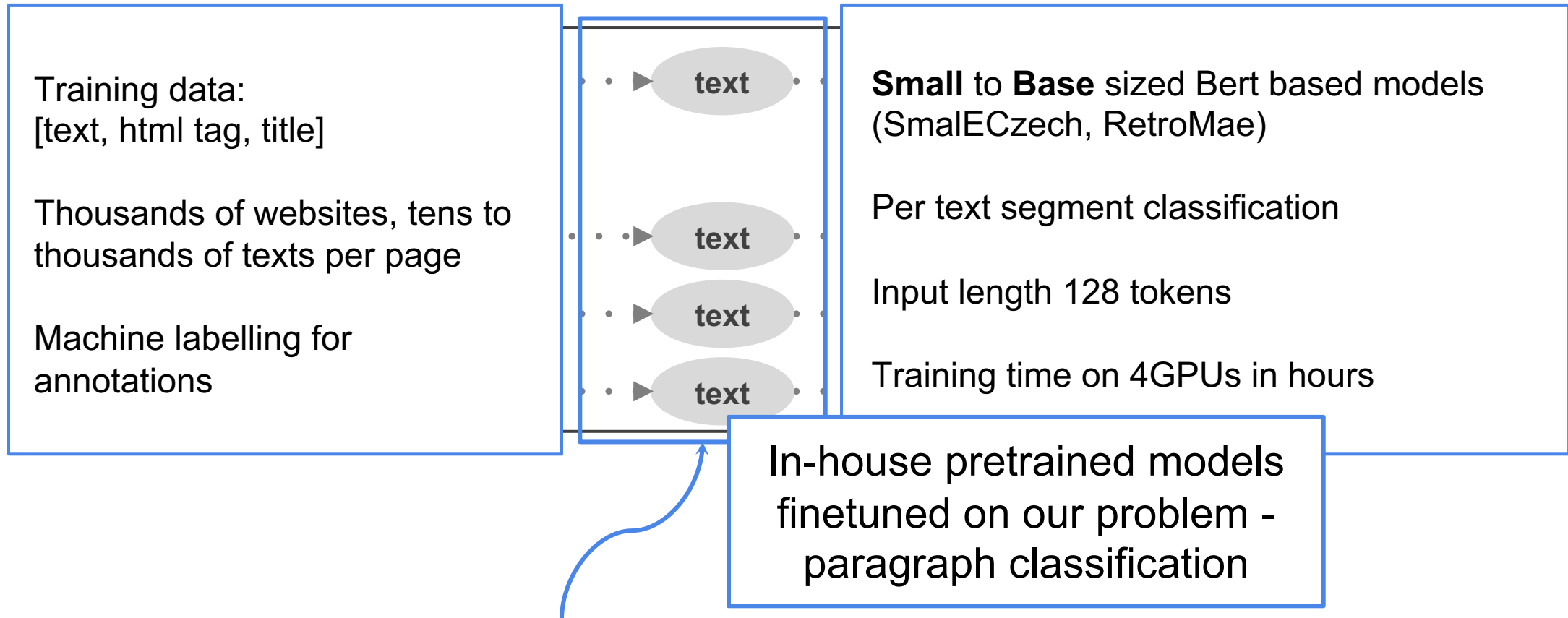
Split



Split - Extract - Vectorize - Evaluate



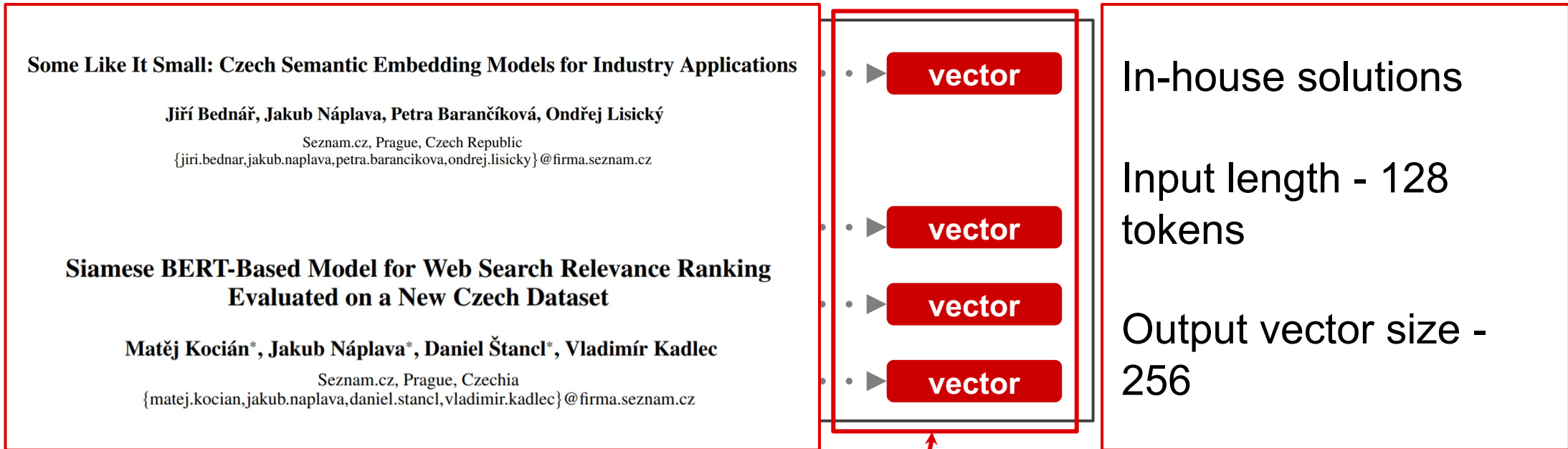
Extract



Split - Extract - Vectorize - Evaluate



Vectorize



Split - Extract - **Vectorize** - Evaluate



Evaluate

Retrieval task:

Retrieve k documents
How many of good
documents did we
retrieved?

rec@k

1 0 1 0 1 1 | 0 0 1 0 0

Relevance task:

Use our text selection in
relevance sorting.
prec@k

1 2 4 3 5 7 | 9 8 10

Evaluation on
downstream
tasks

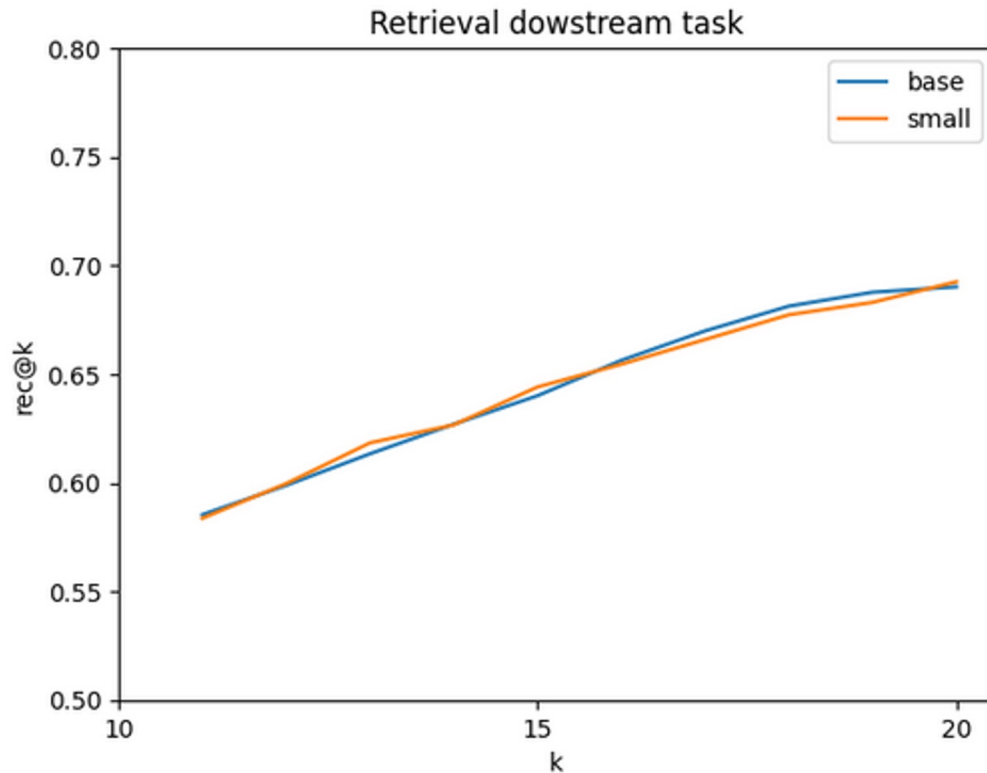
Split - Extract - **Vectorize** - Evaluate



Results



Evaluate Results (very preliminary)



Relevance sorting task		
	small	base
prec@k (mod.)	0.692	0.7

<https://zvukoveizolacnimaterialy.blogspot.com/2016/12/lamino-desky-na-miru.html>

Prodej molitanu na míru Praha Molitanové desky patří k nejprodávanějším druhům pěn na trhu s čalounickým a dalším materiálem. Jedná se nejen o kvalitní desky z molitanové pur pěny, ...

Deskový materiál na dřevěné bázi (překližky, lamino desky, dřevotřísky) a jejich následná úprava, to je Dřevobis, s. Firmám, ale i běžným spotřebitelům jsme schopni. Laminované dřevotřískové desky ...

Dekory lamina a hrany se dají volně kombinovat. Kombinace barev si můžete vyzkoušet v aplikaci dodavatele lamin Egger. Vyberte si z přírodních materiálů jako je borovice, dub, buk, olše, třešeň, u dýhy a lamina je v nabídce ...

....

Small vs Base sized models

Potential to improve



Resume



Self Question-Answering session

Why paragraph extraction approach over text summarization?

Easier, cheaper, in-house (at the moment), the approach do not hallucinate, the approach is `readable`.



Self Question-Answering session

Why paragraph extraction approach over text summarization?

Easier, cheaper, in-house (at the moment), the approach do not hallucinate, the approach is `readable`.

How can we be sure that model can work with separate paragraphs only?

Document is represented by its `html.header.title`, which goes as an input to training and to inference as well.



Self Question-Answering session

Why paragraph extraction approach over text summarization?

Easier, cheaper, in-house (at the moment), the approach do not hallucinate, the approach is `readable`

How can we be sure that model can work with separate paragraphs only?

Document is represented by its `html.header.title`, which goes as an input to training and to inference as well.

How precise does it have to be?

Errors propagate, BUT errors are also suppressed in later search steps: term search, relevance filtering, sorting...



Self Question-Answering session

Why paragraph extraction approach over text summarization?

Easier, cheaper, in-house (at the moment), the approach do not hallucinate, the approach is 'readable'

How can we be sure that model can work with separate paragraphs only?

Document is represented by its html.header.title, which goes as an

How precise does it have to be?

Errors propagate, BUT errors are also suppressed in later search steps: text search, relevance filtering, sorting...

Then, why larger, or even generative models?

Better text selection (smarter than html segment), more applications (direct answers, related questions...), good for wide range of document types.



Future Plans



Future work - research and application

- ✓ Having fun with data: more data, better data, various url types...
- ✓ Having fun with models: small, base, distilled small, longer context, HTML language models
- ✓ Having fun with methodology: sequences instead of paragraphs
- ✓ Having fun with LLM: apply in-house LLM solution (e.g. for better web types coverage)





Martin Bachura

ML Research Team Lead

martin.bachura@firma.seznam.cz

Copyright © 1996–2023 Seznam.cz, a. s.



Prostor pro dotazy

