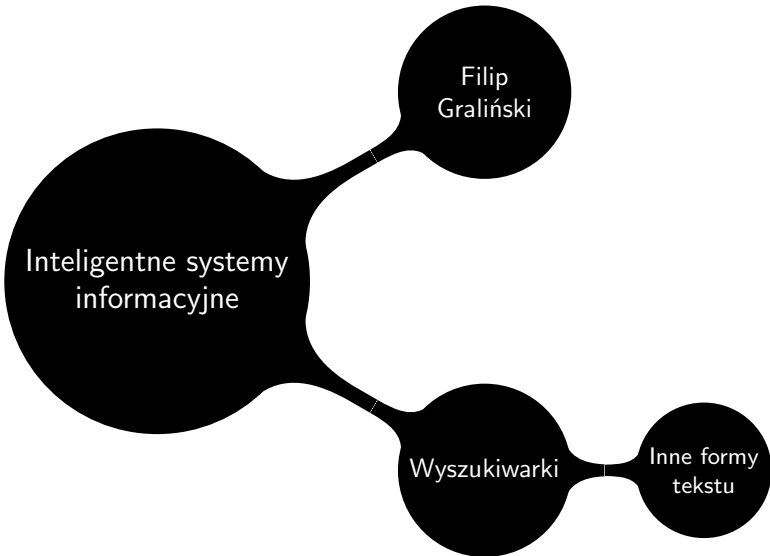


```
graph LR; A[Inteligentne systemy informacyjne] --- B[Filip Graliński]; A --- C[Wyszukiwarki]
```

Inteligentne systemy  
informacyjne

Filip  
Graliński

Wyszukiwarki




Filip  
Graliński

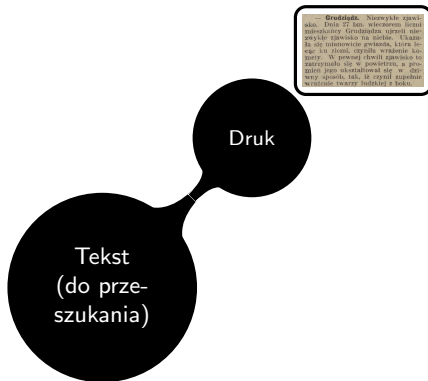
Inteligentne systemy  
informacyjne

Wyszukiwarki

Inne formy  
tekstu



Tekst  
(do prze-  
szukania)





Tekst  
(do przeszukania)

Druk

Pismo  
odręczne

Mowa

— Brudziła. Nowyżło zpaśko. Dula 27 km. wierszem liczył milionowy Grudniaż ugrwił możyło zjawisko na niebie. Ukazała się nitonowicie gwiazda, która legła. Ina ziemia, cyprala, wrodzone kaktusy. W prawej chwili zjawisko to rozprzynało się w powietrzu, a promień jego skłoniłował się w kierunku spowodu, tak, iż czuł się zupełnie wrodzone iwarz. Inicjała z boku.



*Mylorim Mianko: M  
oy obropnie emy.  
Maly i rajunij*

조선글  
한글

— Brudzięda, Nowyżyło zjawisko. Dala 27 km. wierszem liczył mianowicie Grudzięda ugrawi mianowicie zjawisko na niebie. Ukazała się mianowicie gwiazda, która leżała na ziemi, cyprala, widać się kłopoty. W pierwszej chwili zjawisko to rozprysło się w powietrze, a potem jego kształtował się w kierunku spowodu, tak, iż cypral zupełnie widocznie tworzył kształt z boku.



*My brim Mianko: M  
cy obropnie myj  
Wdy i rajunij*

조선글  
한글

— Brudzięda, Nowyżyte zpańsko. Dula 27 kaa, wiczasem licził mianowicie Grudzięda ugawił mianowicie zjawisko na niebie. Ukazało się mianowicie gwiazda, która leżała na ziemi, czerpała, widać się kłóżyła. W prawej chwili zjawisko to zmierzwiło się w powietrze, a potem jego kształtował się w kierunku spowodu, tak, iż czuł się podobnie widać było twardy, białej z boku.

Język obcy

Druk

Tekst  
(do przeszukania)

Texting



Mowa

Pismo odręczne



Myślami Mianko i Mianko  
czy obywatel mianko  
Mianko i rajmuj



## Forma zapytania czy przeszukiwanych treści?

	forma.....	
	zapytania	treści
druk		
pismo odręcz- ne		
mowa		
język obcy		
<i>texting</i>		

## Forma zapytania czy przeszukiwanych treści?

	forma.....	
	zapytania	treści
druk	—	przeszukiwanie historycznych materiałów
pismo odręczne		
mowa		
język obcy		
<i>texting</i>		

## Forma zapytania czy przeszukiwanych treści?

	forma. . . . .	
	zapytania	treści
druk	—	przeszukiwanie historycznych materiałów
pismo odręczne	interfejs urządzeń mobilnych (?)	przeszukiwanie rękopisów (??)
mowa		
język obcy		
<i>texting</i>		

## Forma zapytania czy przeszukiwanych treści?

	forma. . . . .	
	zapytania	treści
druk	—	przeszukiwanie historycznych materiałów
pismo odręczne	interfejs urządzeń mobilnych (?)	przeszukiwanie rękopisów (??)
mowa	<i>voice search</i>	<i>spoken content retrieval</i>
język obcy		
<i>texting</i>		

## Forma zapytania czy przeszukiwanych treści?

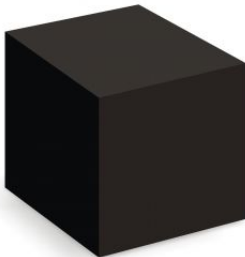
	forma.....	
	zapytania	treści
druk	—	przeszukiwanie historycznych materiałów
pismo odręczne	interfejs urządzeń mobilnych (?)	przeszukiwanie rękopisów (??)
mowa	<i>voice search</i>	<i>spoken content retrieval</i>
język obcy	—	<i>cross-language information retrieval</i>
<i>texting</i>		

## Forma zapytania czy przeszukiwanych treści?

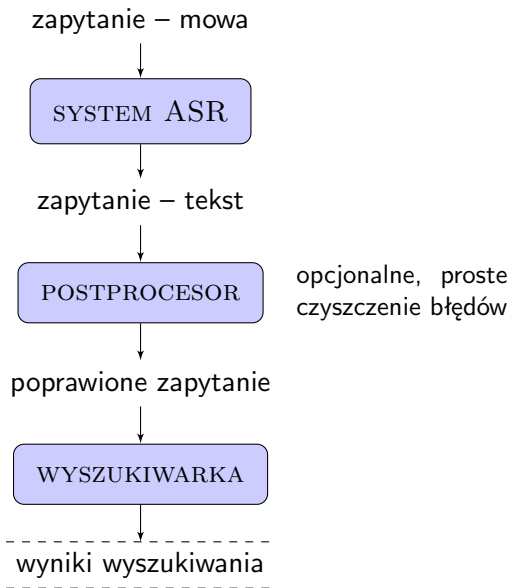
	forma. . . . .	
	zapytania	treści
druk	—	przeszukiwanie historycznych materiałów
pismo odręczne	interfejs urządzeń mobilnych (?)	przeszukiwanie rękopisów (??)
mowa	<i>voice search</i>	<i>spoken content retrieval</i>
język obcy	—	<i>cross-language information retrieval</i>
<i>texting</i>	interfejs urządzeń mobilnych	przeszukiwanie Twittera, chatów, SMS-ów etc.

# Jak wyszukiwać w innych formach tekstu?

Rozwiązanie 1. Najprościej traktować OCR/HWR/ASR/MT jako czarną skrzynkę, która zwraca zwykły tekst.

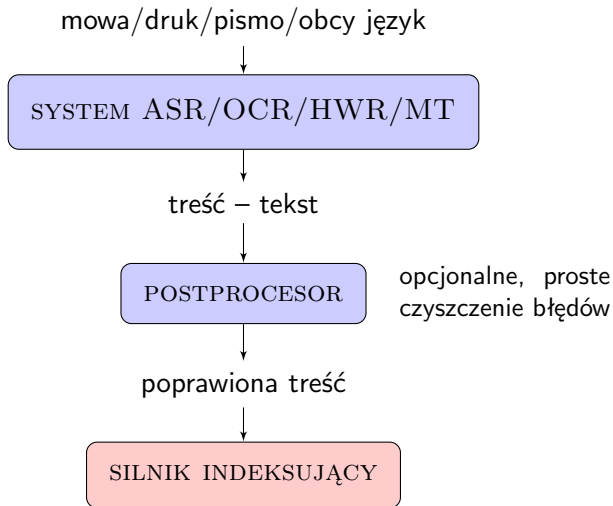


## Przez czarną skrzynkę przepuszczamy zapytanie. . .

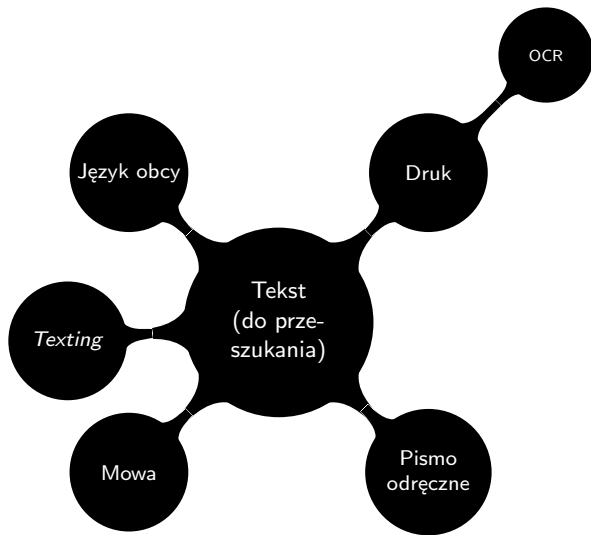




... albo indeksowane treści



## Przyjrzyjmy się bliżej czarnym skrzynkom...



# OCR = Optical Character Recognition

Trudności:

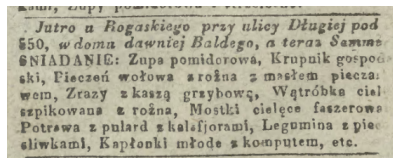
*Jutro u Bogaskiego przy ulicy Długiej pod  
550, w domu dawniej Baldego, a teraz Sammt  
ŚNIADANIE: Zupa pomidorowa, Krupnik gospo  
ski, Pieczeń wołowa z rożną z masłem piecza  
wem, Zrazy z kaszą grybową, Wątróbka ciel  
szpikowana z rożną, Mostki cielęce faszerowa  
Potrawa z pulard z kalafjorami, Legumina z pie  
śliwkami, Kapłonki młode z komputem, etc.*

*Jutro u Bogaskiego przy ulicy Długiej pod  
550, w domu dawniej B aide go, a teraz Sammt  
ŚNIADANIE: Zupa pomidorowa, Krupnik  
gospoó ski, Pieczeń wołowa z rożną 5 msstetn  
piecza; wetn, Zrazy z kaszą grybowy,  
Wątróbka ciel i szpikowana Z rożna, Mostki  
cielęce faszerowp Potrawa z pulard %  
kalafjorami, Legnmina z pie śliwkami, Kapłonki  
młode z komputera, etc.*

# OCR = Optical Character Recognition

Trudności:

- ▶ zniekształcone strony, kleksy, pleśń, dżemik itd.



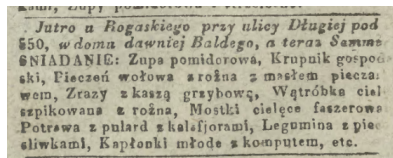
Jutro u Bogaskiego przy ulicy Długiej pod  
550, w domu dawniej Baldego, a teraz Sammt  
ŚNIADANIE: Zupa pomidorowa, Krupnik gospo  
ski, Pieczeń wołowa z rożną z masłem piecza  
wem, Zrazy z kaszą grybową, Wątróbka ciel  
szpikowana z rożną, Mostki cielęce faszerowa  
Potrawa z pulard z kalafjorami, Legumina z pie  
śliwkami, Kapłonki młode z komputem, etc.

*Jutro u Bogaskiego przy ulicy Długiej pod  
550, w domu dawniej Baldego, a teraz Sammt  
ŚNIADANIE: Zupa pomidorowa, Krupnik  
gospoó ski, Pieczeń wołowa z rożną 5 msłtetn  
piecza; wetn, Zrazy z kaszą grybowy,  
Wątróbka ciel i szpikowana Z rożna, Mostki  
cielęce faszerowp Potrawa z pulard %  
kalafjorami, Legnmina z pie śliwkami, Kapłonki  
młode z komputera, etc.*

# OCR = Optical Character Recognition

## Trudności:

- ▶ zniekształcone strony, kleksy, pleśń, dżemik itd.
- ▶ rozpoznanie struktury strony (kolumny, ilustracje itd.)



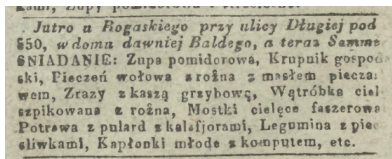
Jutro u Bogaskiego przy ulicy Długiej pod  
550, w domu dawniej Baldego, a teraz Sammt  
ŚNIADANIE: Zupa pomidorowa, Krupnik gospo  
ski, Pieczeń wołowa z rożna z masłem piecza  
wem, Zrazy z kaszą grybową, Wątróbka ciel  
szpikowana z rożna, Mostki cielęce faszerowa  
Potrawa z pulard z kalafjorami, Legumina z pie  
śliwkami, Kapłonki młode z komputem, etc.

*Jutro u Bogaskiego przy ulicy Długiej pod  
550, w domu dawniej B aide go, a teraz Sammt  
ŚNIADANIE: Zupa pomidorowa, Krupnik  
gospoó ski, Pieczeń wołowa z rolna 5 msstetn  
piecza; wetn, Zrazy z kaszą grybowy,  
Wątróbka ciel i szpikowana Z rożna, Mostki  
cielęce faszerowp Potrawa z pulard %  
kalafjorami, Legnmina z pie śliwkami, Kapłonki  
młode z komputera, etc.*

# OCR = Optical Character Recognition

## Trudności:

- ▶ zniekształcone strony, kleksy, pleśń, dżemik itd.
- ▶ rozpoznanie struktury strony (kolumny, ilustracje itd.)
- ▶ zmieniające się konwencje typograficzne, kroje czcionek itd.
- ▶ zmieniające się konwencje ortograficzne

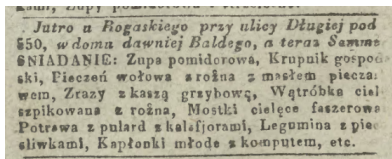


*Jutro u Bogaskiego przy ulicy Długiej pod  
550, w domu dawniej B aide go, a teraz Sammt  
ŚNIADANIE: Zupa pomidorowa, Krupnik  
gospoó ski, Pieczeń wołowa z rożna 5 msstetn  
piecza; wetn, Zrazy z kaszą grybowy,  
Wątróbka ciel i szpikowana Z rożna, Mostki  
cielęce faszerowp Potrawa z pulard %  
kalafjorami, Legnmina z pie śliwkami, Kapłonki  
młode z komputera, etc.*

# OCR = Optical Character Recognition

## Trudności:

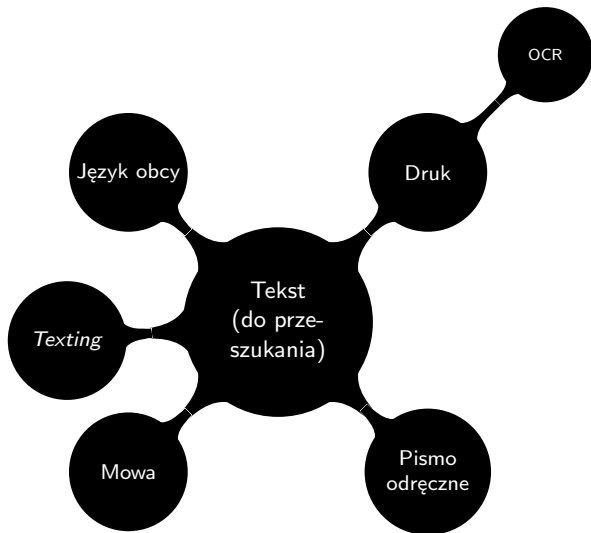
- ▶ zniekształcone strony, kleksy, pleśń, dżemik itd.
- ▶ rozpoznanie struktury strony (kolumny, ilustracje itd.)
- ▶ zmieniające się konwencje typograficzne, kroje czcionek itd.
- ▶ zmieniające się konwencje ortograficzne



*Jutro a Bogaskiego przy ulicy Długiej pod 550, w domu dawniej Baldego, a teraz Sammt*  
*ŚNIADANIE: Zupa pomidorowa, Krupnik gospoó ski, Pieczeń wołowa z rolna 5 msstetn piecza; wetn, Zrazy z kaszą grzybowy, Wątróbka ciel i szpikowana Z rożna, Mostki cielęce faszerowp Potrawa z pulard % kalafjorami, Legnmina z pie śliwkami, Kapłonki młode z komputera, etc.*

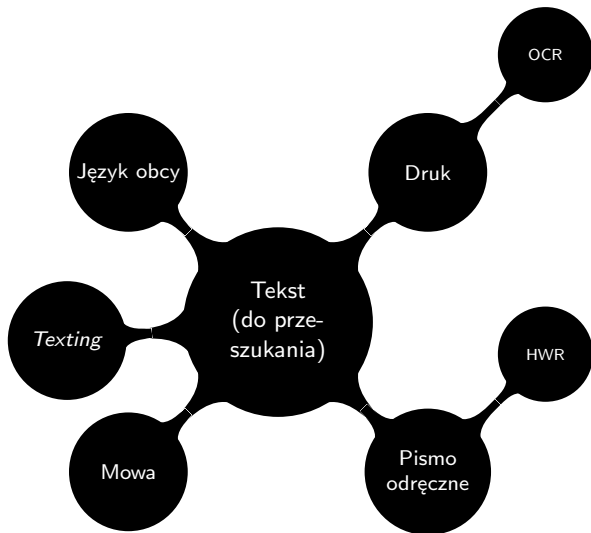
Przydatne narzędzia: unpaper (usuwa niektóre zniekształcenia), tesseract (OCR), ocrdju (OCR-owanie plików DjVu)

## Przyjrzyjmy się bliżej czarnym skrzynkom...

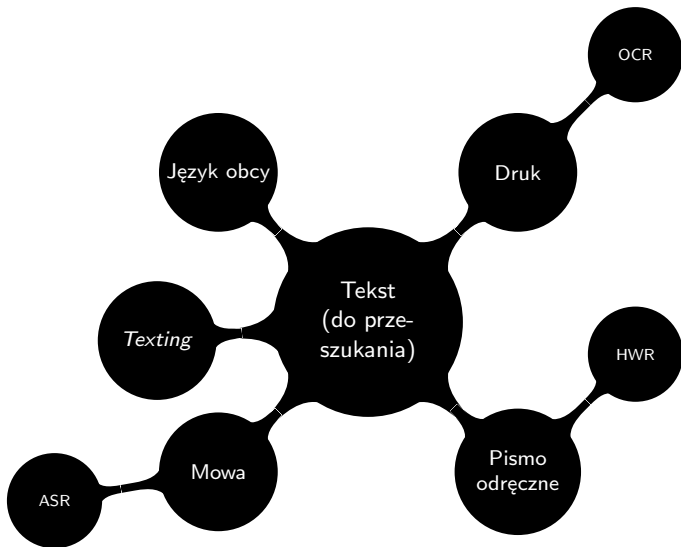




## Przyjrzyjmy się bliżej czarnym skrzynkom...



## Przyjrzyjmy się bliżej czarnym skrzynkom...



# Rozpoznawanie mowy (ASR, *Automatic Speech Recognition*)

Trudności:

- ▶ różnice w wymowie (indywidualne, dialekty, płeć itd.)

# Rozpoznawanie mowy (ASR, *Automatic Speech Recognition*)

Trudności:

- ▶ różnice w wymowie (indywidualne, dialekty, płeć itd.)
- ▶ szumy tła i urządzenia, niedoskonałości mikrofonu itd.

# Rozpoznawanie mowy (ASR, *Automatic Speech Recognition*)

Trudności:

- ▶ różnice w wymowie (indywidualne, dialekty, płeć itd.)
- ▶ szumy tła i urządzenia, niedoskonałości mikrofonu itd.
- ▶ nieregularności pisowni (nazwy własne obcojęzyczne!)

# Rozpoznawanie mowy (ASR, *Automatic Speech Recognition*)

Trudności:

- ▶ różnice w wymowie (indywidualne, dialekty, płeć itd.)
- ▶ szумы tła i urządzenia, niedoskonałości mikrofonu itd.
- ▶ nieregularności pisowni (nazwy własne obcojęzyczne!)
- ▶ wymaga dużej mocy obliczeniowej

# Rozpoznawanie mowy (ASR, *Automatic Speech Recognition*)

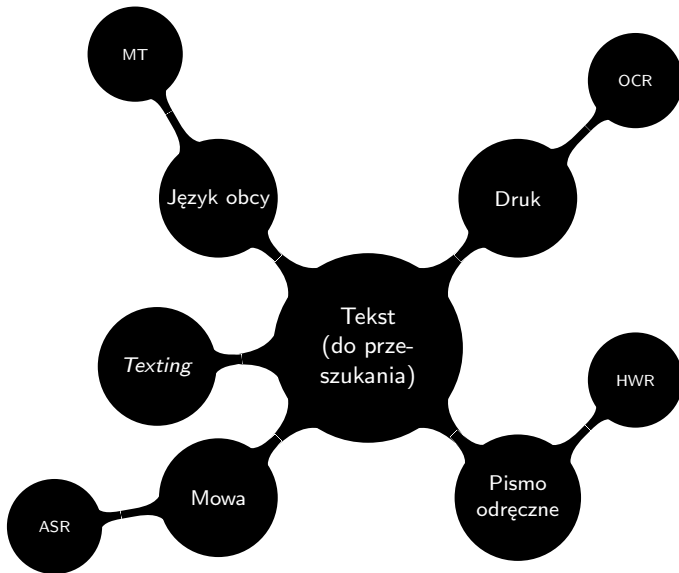
## Trudności:

- ▶ różnice w wymowie (indywidualne, dialekty, płeć itd.)
- ▶ szумы tła i urządzenia, niedoskonałości mikrofonu itd.
- ▶ nieregularności pisowni (nazwy własne obcojęzyczne!)
- ▶ wymaga dużej mocy obliczeniowej

## Przydatne narzędzia:

- ▶ do tworzenia systemów ASR: CMU Sphinx, Kaldi
- ▶ Google Speech API (pokątnie)

## Przyjrzyjmy się bliżej czarnym skrzynkom...





# Tłumaczenie maszynowe (MT, *machine translation*)

Dlaczego trudne?

- ▶ niejednoznaczności w tekście oryginalnym
  - ▶ homonimy: *powieść, Łódź, nie, lecz, albo*
  - ▶ co do części mowy: *komputerowi, pożółkły, last*
  - ▶ składniowa: *I saw a girl with a telescope*
  - ▶ semantyczna: *mysz, kanał, jack, table*
- ▶ niejednoznaczności między językami
  - ▶ *palec = finger czy toe (digit?)*
  - ▶ *rzeka = fleuve czy riviere?*
  - ▶ *siostra = meimei czy jiejie (jiemei?)*
  - ▶ *zobaczyłem dziewczynę = I saw a girl czy I saw the girl?*

## A texting?

- ▶ przywracanie diakrytyków (*diacritic restoration*)
- ▶ specyficzne skróty (*LOL, 4 U, 3maj się*)
- ▶ specyficzne konwencje (np. hasztagi)

## Zajrzyjmy do wnętrza czarnych skrzynek. . .

- ▶ ASR i MT – model zaszumionego kanału (*noisy channel model*)

## Zajrzyjmy do wnętrza czarnych skrzynek. . .

- ▶ ASR i MT – model zaszumionego kanału (*noisy channel model*)
  - ▶ mowa to zniekształcona forma tekstu

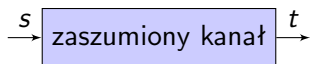
## Zajrzyjmy do wnętrza czarnych skrzynek. . .

- ▶ ASR i MT – model zaszumionego kanału (*noisy channel model*)
  - ▶ mowa to zniekształcona forma tekstu
  - ▶ tekst np. koreański jest tylko dziwnie zniekształconą formą tekstu polskiego

## Zajrzyjmy do wnętrza czarnych skrzynek. . .

- ▶ ASR i MT – model zaszumionego kanału (*noisy channel model*)
  - ▶ mowa to zniekształcona forma tekstu
  - ▶ tekst np. koreański jest tylko dziwnie zniekształconą formą tekstu polskiego
- ▶ HWR i OCR – model zaszumionego kanału ma pewne zastosowanie, ale istotne jest przetwarzanie w 2 wymiarach

## Model zaszumionego kanału



tłumaczenia / mowa

$$P(t) = P(s)P(t|s)$$

ale nas interesuje  $P(s|t)$ !

## Piękno i wierność

$$P(e|f) = \frac{P(f|e)P(e)}{P(f)}$$



## Piękno i wierność

$$P(e|f) = \frac{P(f|e)P(e)}{P(f)}$$

- ▶ Jak dobre jest tłumaczenie?

$$P(e|f) \propto P(f|e) * P(e)$$

- ▶ Wierność

# Piękno i wierność

$$P(e|f) = \frac{P(f|e)P(e)}{P(f)}$$

- ▶ Jak dobre jest tłumaczenie?

$$P(e|f) \propto P(f|e) * P(e)$$

- ▶ Wierność
- ▶ Piękno

$$\tilde{e} = \arg \max_{e \in e^*} P(f|e)P(e)$$

## Modele języka

- ▶  $p(f|e)$  będzie szacowane za pomocą **modelu tłumaczenia** (*translation model*)
- ▶  $p(e)$  będzie szacowane za pomocą modelu języka (docelowego) (*target language model*)

## Modele języka

- ▶  $p(f|e)$  będzie szacowane za pomocą modelu tłumaczenia (*translation model*)
- ▶  $p(e)$  będzie szacowane za pomocą **modelu języka (docelowego)** (*target language model*)

## Modele języka

- ▶  $p(f|e)$  będzie szacowane za pomocą modelu tłumaczenia (*translation model*)
- ▶  $p(e)$  będzie szacowane za pomocą modelu języka (docelowego) (*target language model*)

**Model języka** przypisuje prawdopodobieństwa napisom.

Jeśli  $M$  ma być modelem języka polskiego, oczekiwalibyśmy, że dla napisów:

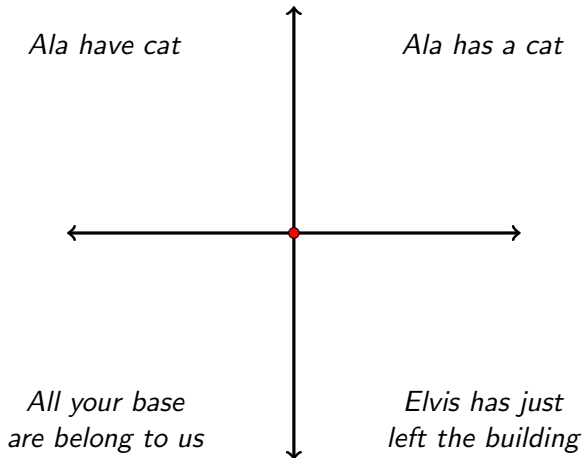
- $z_1$  *W tym stanie rzeczy pan Ignacy coraz częściej myślał o Wokulskim.*
- $z_2$  *Po wypełniony zbiornik pełny i należne kwotę, usłyszała w attendant*
- $z_3$  *xxxxyżzzzzit backspace hoooooooooop x y z*

zachodzić będzie:

$$P(z_1|M) > P(z_2|M) > P(z_3|M)$$

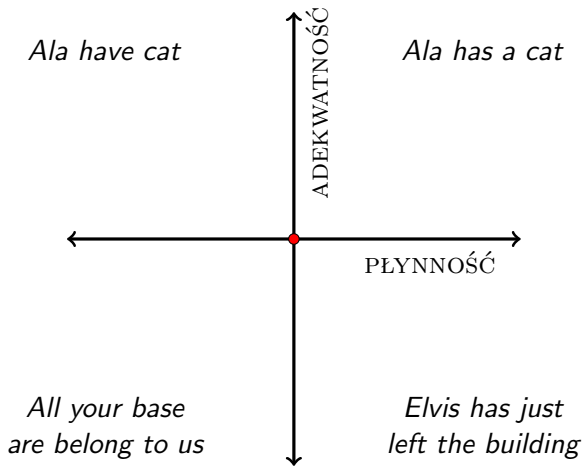
## Dwa wymiary tłumaczenia

*Ala ma kota = ...*

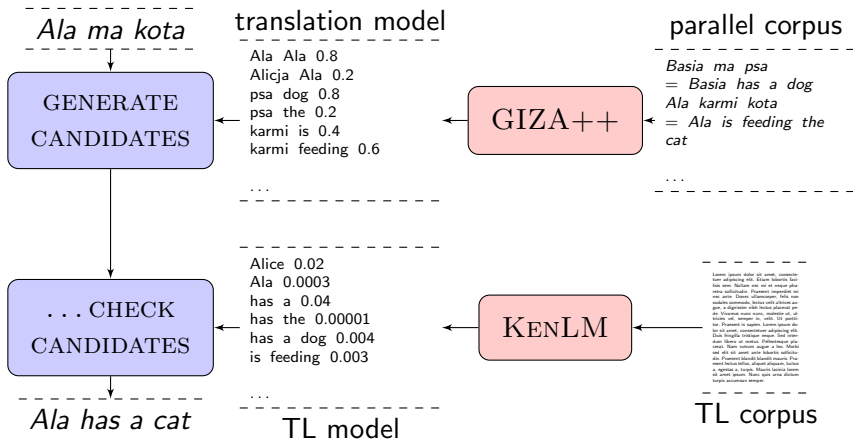


# Dwa wymiary tłumaczenia

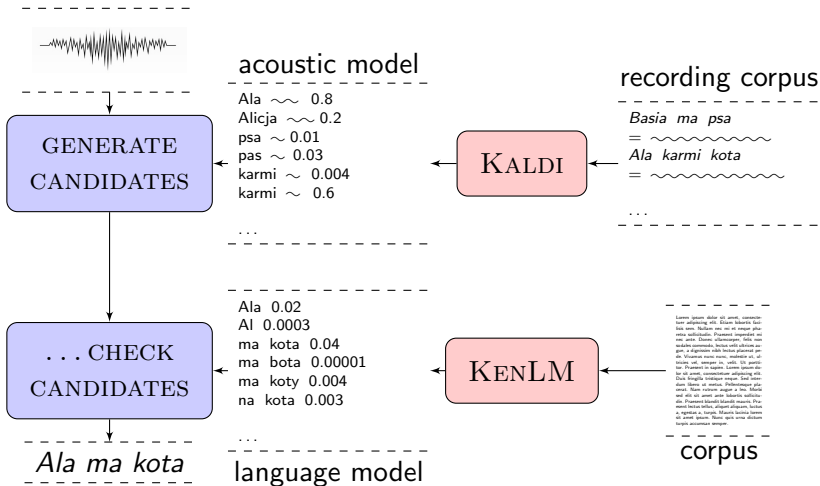
*Ala ma kota = ...*



# Statystyczne tłumaczenia maszynowe







## Jak wyszukiwać w innych formach tekstu?

Rozwiązanie 1. Najprościej traktować OCR/HWR/ASR/MT jako czarną skrzynkę, która zwraca zwykły tekst.

Rozwiązanie 2. Brać pod uwagę *n-best list*, nie tylko pierwszą odpowiedź.

## lista $n$ najlepszych ( $n$ -best list)

Przykładowa 4-best list (wyjście z ASR-a):

	<b>możliwa interpretacja</b>	<b>prawdopodobieństwo</b>
1.	<i>wojewódzki narodowy o powołanie komisji ochrony przyrody krajowej</i>	0.5219596
2.	<i>rejowa wojewódzki narodowy o powołanie komisji ochrony przyrody krajowej</i>	
3.	<i>prawa wojewódzki narodowy o powołanie komisji ochrony przyrody krajowej</i>	
4.	<i>wojewódzki narodowy o powołanie komisji ochrony przyrody krajowe</i>	

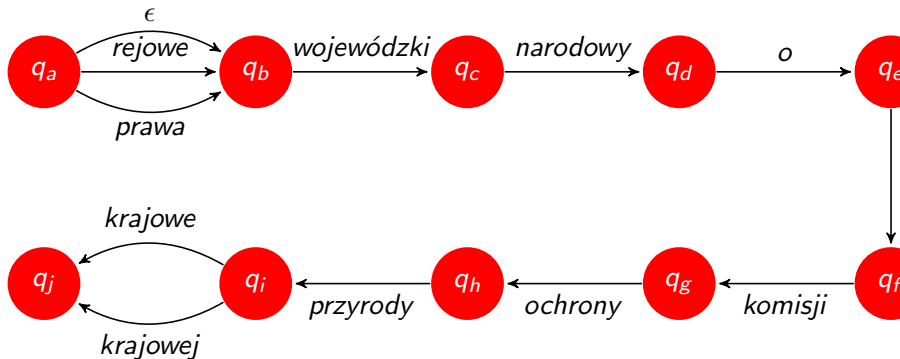
# Jak wyszukiwać w innych formach tekstu?

Rozwiązanie 1. Najprościej traktować OCR/HWR/ASR/MT jako czarną skrzynkę, która zwraca zwykły tekst.

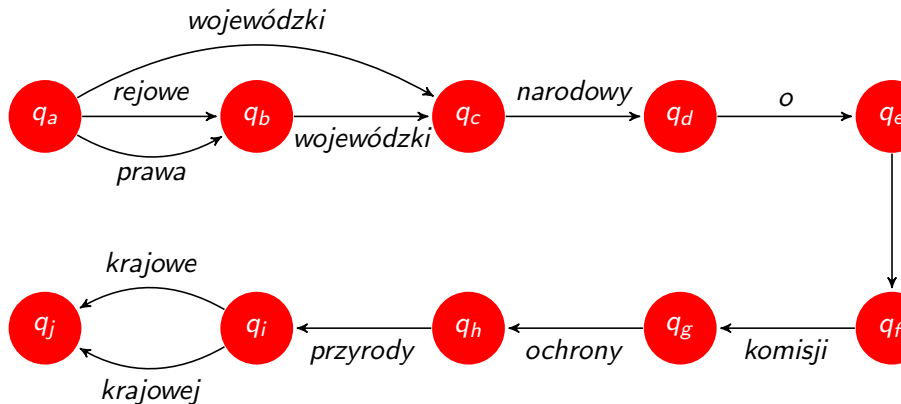
Rozwiązanie 2. Brać pod uwagę *n-best list*, nie tylko pierwszą odpowiedź.

Rozwiązanie 3. ... albo bardziej zaawansowane sposoby reprezentacji wyjścia – „**kiełbaski**” albo **kraty**.

# „Kiełbaski” (word confusion networks)



# Kraty (lattices)



## Jak wyszukiwać w innych formach tekstu?

Rozwiązanie 1. Najprościej traktować OCR/HWR/ASR/MT jako czarną skrzynkę, która zwraca zwykły tekst.

Rozwiązanie 2. Brać pod uwagę *n-best list*, nie tylko pierwszą odpowiedź.

Rozwiązanie 3. ... albo bardziej zaawansowane sposoby reprezentacji wyjścia – „kiełbaski” albo kraty.

Rozwiązanie 4. **Ściśle sprząc modele OCR/HWR/ASR/MT i wyszukiwania**

## Przykład: wyszukiwanie lokalne (*local search*)

$$\hat{L} = \arg \max_L P(L)P(C|L)P(Q|L)P(Q|X)$$

- ▶  $L$  – miejsce
- ▶  $C$  – gdzie przebywa użytkownik
- ▶  $Q$  – zapytanie (tekst)
- ▶  $X$  – sygnał dźwiękowy



## Jak wyszukiwać w innych formach tekstu?

Rozwiązanie 1. Najprościej traktować OCR/HWR/ASR/MT jako czarną skrzynkę, która zwraca zwykły tekst.

Rozwiązanie 2. Brać pod uwagę *n-best list*, nie tylko pierwszą odpowiedź.

Rozwiązanie 3. ... albo bardziej zaawansowane sposoby reprezentacji wyjścia – „kiełbaski” albo kraty.

Rozwiązanie 4. Ściśle sprząc modele OCR/HWR/ASR/MT i wyszukiwania...

Rozwiązanie 5. ... i trenować wszystkie moduły **dyskryminacyjnie**.

Trenowanie dyskryminacyjne – na przykład ASR jest trenowany pod kątem jak najlepszych wyników wyszukiwarki, której jest częścią, nie zmniejszenia swojej stopy błędów.