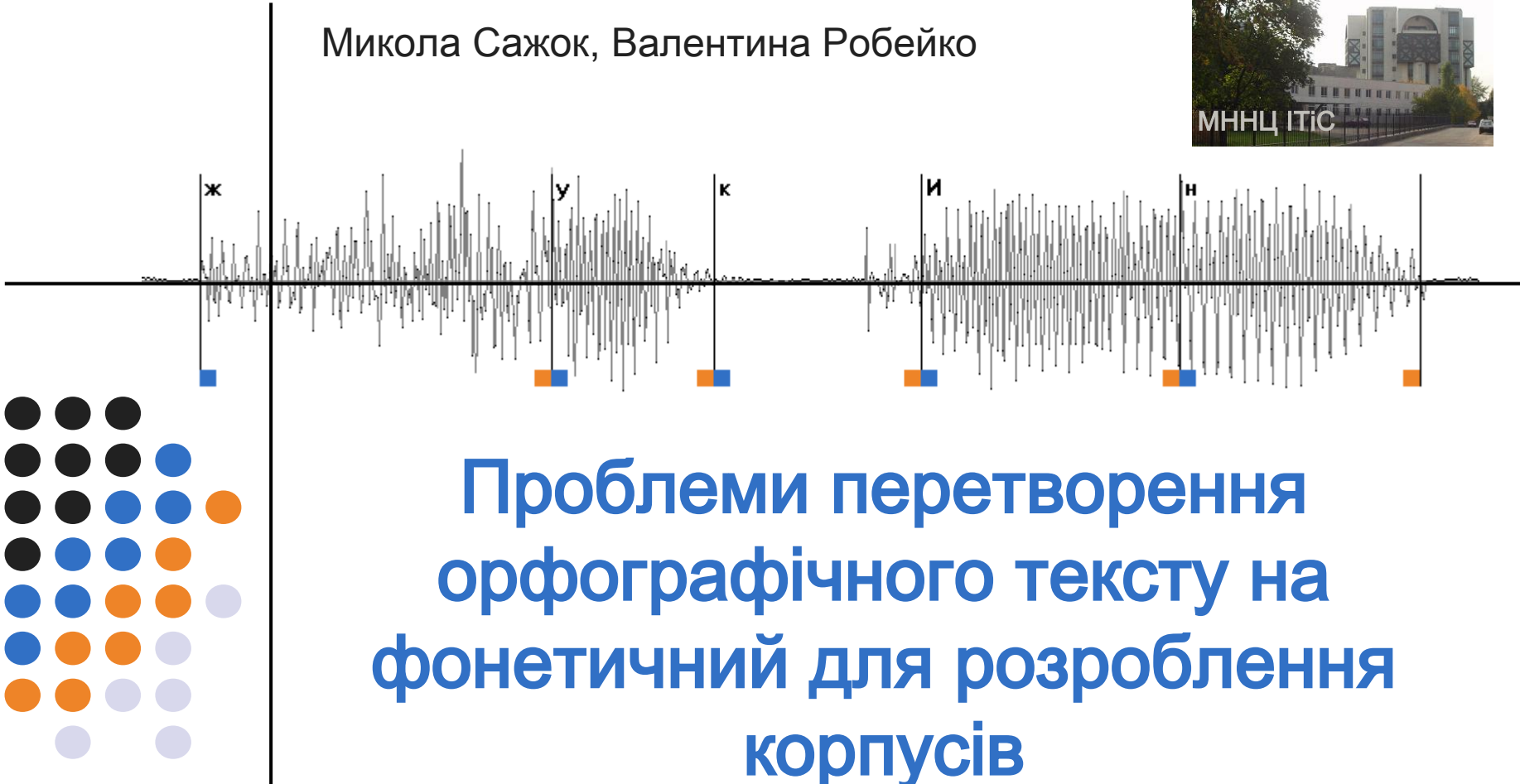


Микола Сажок, Валентина Робейко



# Проблеми перетворення орфографічного тексту на фонетичний для розроблення корпусів

Жукин'2010



## Зміст

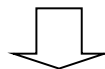
---

- Загальна структура перетворювача графем на фонемі
- Перетворювач
- Розбиття на синтагми
  - Неоднозначність функції токена
- Від графеми до фонемі
  - Гомографія
- Експериментальні результати
- Висновки



## Загальна структура перетворювача графем на фонемні

Мовленнєва інформаційна технологія



Автоматичне перетворення  
орфографічного тексту на фонемний



Гіпотетичний вплив сусідніх  
звуків (слів) у потоці злитого  
мовлення на рівні фонем



Індивідуальність вимови  
кооперативу дикторів



## Загальна структура перетворювача графем на фонемі

Урахування меж між морфемами

Урахування наголосів

Транскриптор

Урахування фонетичних змін на стиках слів у процесі мовлення

Контекстно-залежні правила перетворень між послідовностями символів



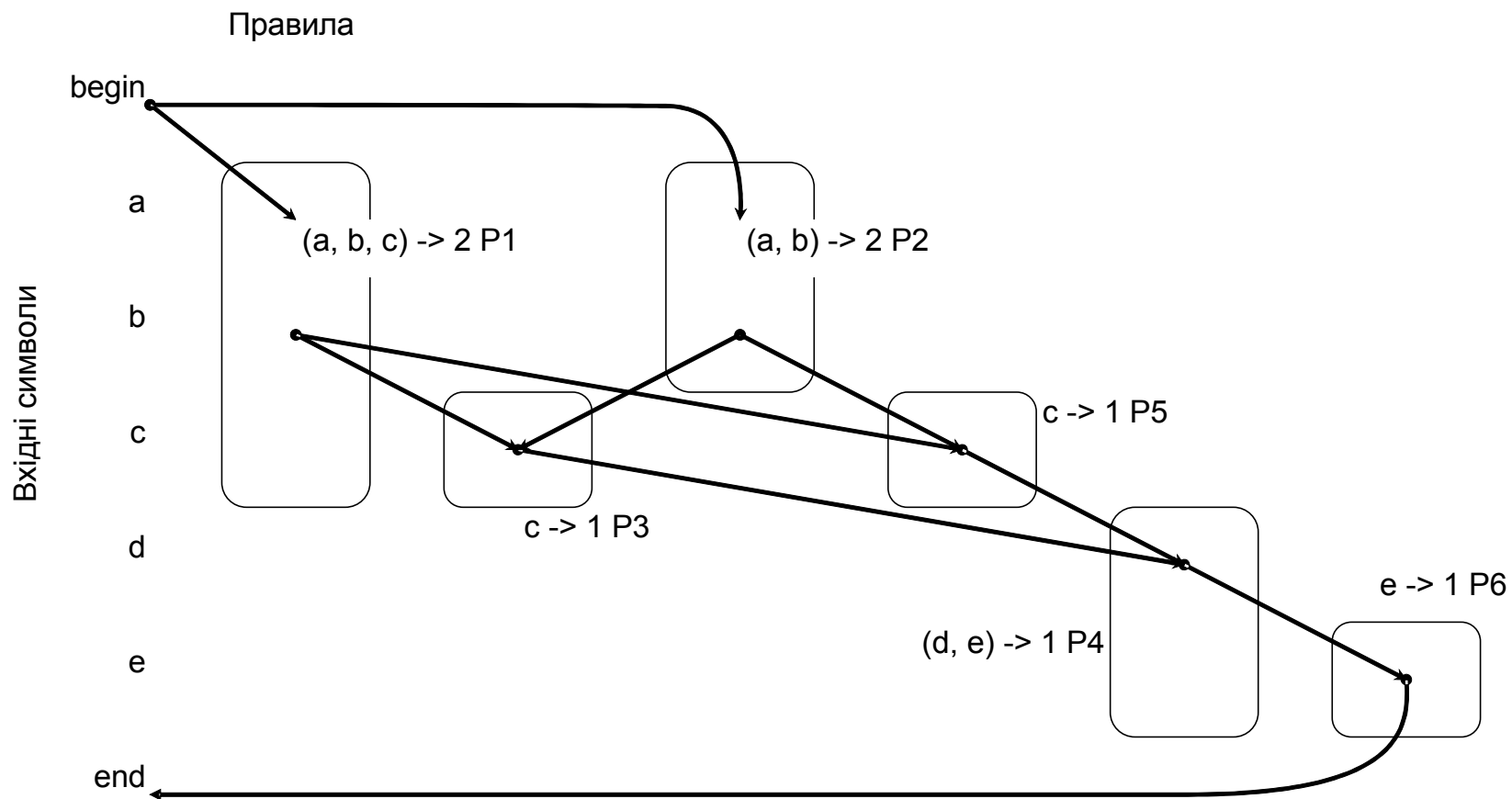
## Стандартні підходи

---

1. Таблиця замін
  - величезна кількість елементів
  - проблемно реалізувати багатозначність
1. Програмні модулі
  - потрібна перекомпіляція для внесення змін
  - нерегуляризованість



# Мультиоператор





## Мультиоператор

---

Параметри:

1. Найдовша довжина вхідної підпослідовності
2. Імовірності (взаємного) застосування кожного з елементарних операторів



## Мультиоператор

---

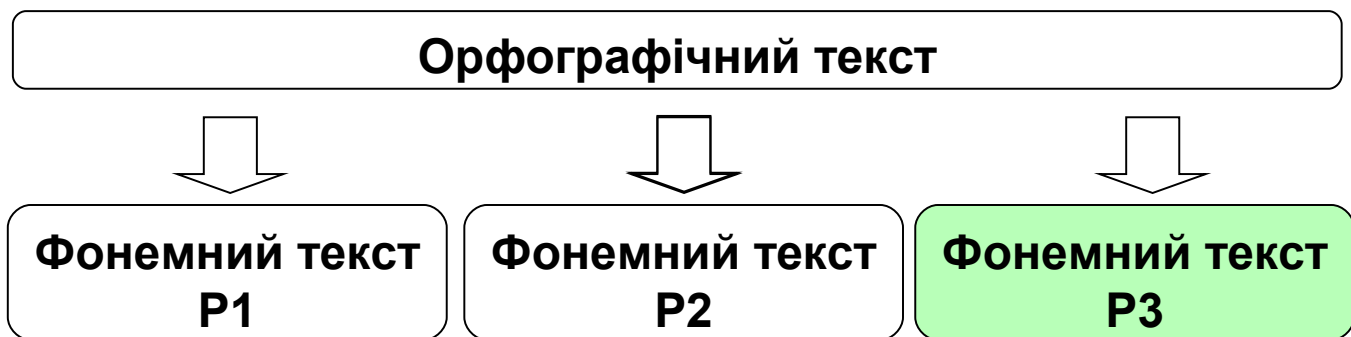
Особливості реалізації:

1. Відома вхідна послідовність повністю
2. Послідовність надходить по відліку
3. Затримка на найдовшу довжину вхідної підпослідовності

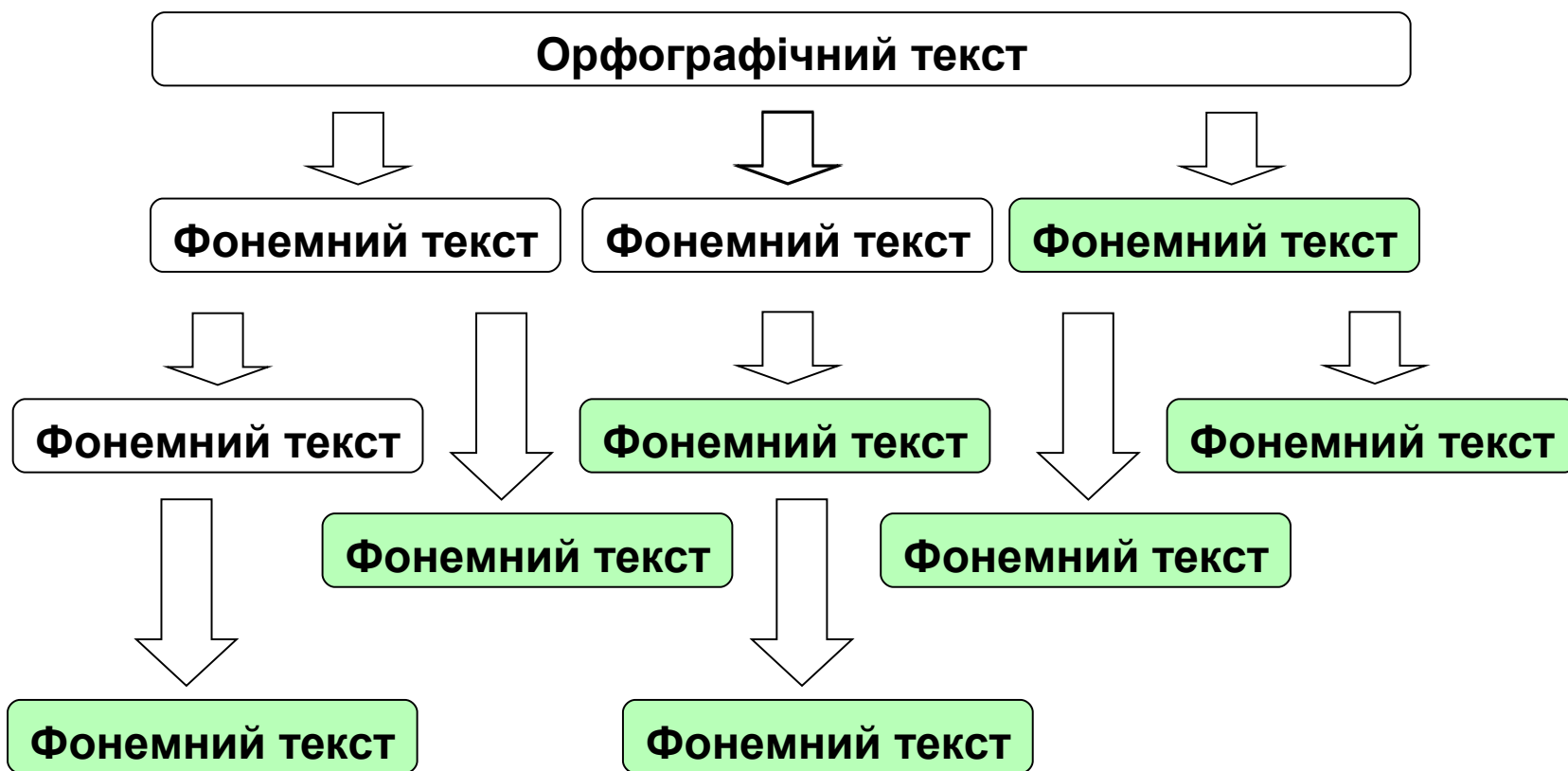




## Варіанти транскрибування



## Варіанти транскрибування – каскад



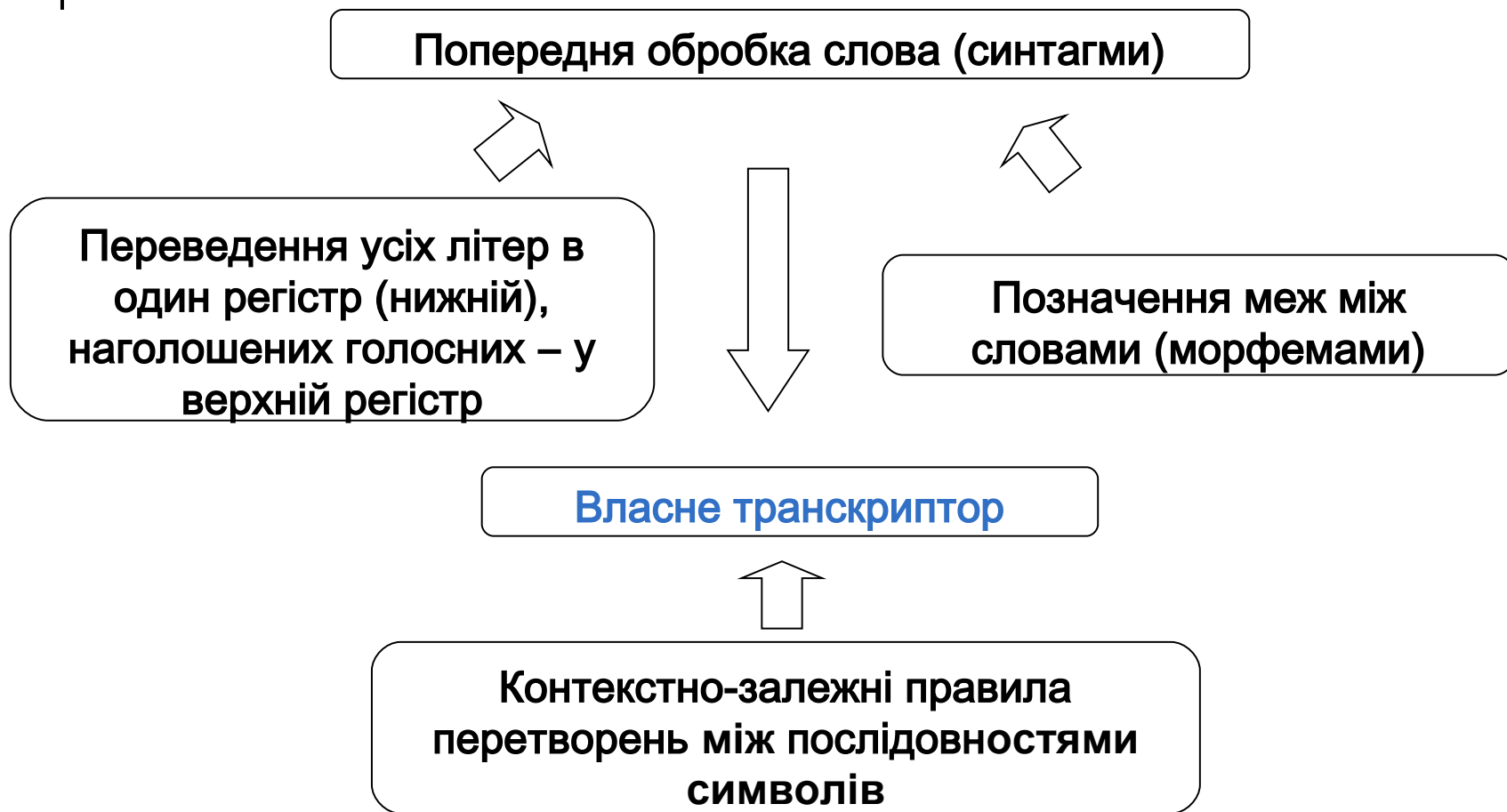


## Варіанти транскрибування: зворотнє





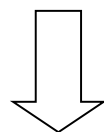
## Транскриптор. Принцип роботи





## Приклади правил побудови транскриптора

Вхідна послідовність графем	Вихідна послідовність фонем	Ширина кроку аналізу	Пояснення
[зсц] [жшч]	[жшч]	1	<b>з, с, ц</b> перед <b>ж, ш, ч</b> переходять відповідно у <b>ж, ш, ч</b>
т[дтзснц] [іієюяЄЮЯь]	т'	1	<b>т</b> перед м'якими <b>д, т, з, с, н, ц</b> пом'якшується
с т [лн]	с	2	<b>т</b> між <b>с</b> та <b>л</b> або <b>н</b> випадає



Близько 30 подібних правил

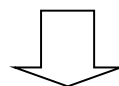
Літня школа-семінар “Усномовні технології та проблеми створення корпусів” Житомир’2010



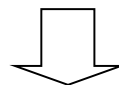
## Варіанти транскрибування

---

**Загальний варіант  
транскрибування на основі  
літературної вимови**



**Розмовні варіанти**

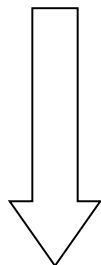


**Індивідуалізовані транскрипції  
для груп дикторів**



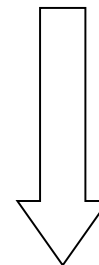
## Позиційні та комбінаторні зміни звуків у потоці мовлення

Позиційні зміни звуків  
у потоці мовлення



боротьба → б а р а д' б А  
          або б р а д' б А  
безпекою → б е с п Е к о у  
доброю → д О б р о й

Комбінаторні зміни  
звуків у потоці  
мовлення



книжка → к н И ш к а  
квітка → к' в' І т к а  
аеропорт → а р о п О р т  
чесний → ч Е с т н и й



## Позиційні зміни звуків у потоці мовлення

- окрім редукції ненаголошених **е, и та о** до **е<sup>и</sup>, и<sup>е</sup>, о<sup>у</sup>**, також ослаблена вимова **о** як **а** в ненаголошеній позиції, рідше трапляється редукція ненаголошених голосних до повного зникнення (**тепер** → **т и п Е р**, **зозуля** → **з у з У л' а**, **боротьба** → **б а р а д' б А** або **б р а д' б А**);
- оглушення дзвінких приголосних перед паузою (**брі́д** → **б р' і т**, **зараз** → **з А р а с**);
- редукція у термінальних частинах слів у процесі мовлення – зникнення приголосного звука в закінченнях **-ого, -их, -ий, -іх, -ій, -ії, -ої, -еї, ою, -єю, -ити** та подібних (**коротший** → **к о р О ч ш и**, **синіх** → **с И н' і**, **безпекою** → **б е с п Е к о у**); зникнення кінцевого голосного звука в закінченнях **-ою, -єю, -єю** та подібних (**доброю** → **д О б р о й**, **землею** → **з е м л Е й**) та ін.





## Комбінаторні зміни звуків у потоці мовлення

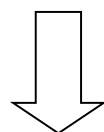
- повна регресивна асиміляція за глухістю у сполучі „дзвінкий+глухий” на межі будь-яких морфем у слові та на межі слів (без причини → б е с п р и ч И н и, розсунути → р о с с У н у т и, книжка → к н И ш к а, сядьте → с' А т' т е);
- асиміляція за м'якістю свистячих та шиплячих приголосних, губних та задньоязикових приголосних (злі → з' л' І, шлях → ш' л' А х, квітка → к' в' І т к а);
- вимова подовжених приголосних звуків як звичайного неподовженого звука, вимова двох голосних як одного звука (віддати → в' і д А т и, знання → з н а н' А, зоопарк → з о п А р к, аеропорт → а р о п О р т);
- неповне спрощення в групах приголосних, його відсутність (чесний → ч Е с т н и й) та ін.

Літня школа-семінар “Усномовні технології та проблеми створення корпусів” Житомир’2010



## Приклади індивідуалізованих словників

### Різновиди транскрипцій словника



<b>Літ. транскрипція</b>	<b>Інд.словник</b>	<b>devocal</b>	<b>end_cons</b>	<b>a</b>	<b>devocal_a</b>
р об И в		р об И ф	р об И	р аб И в	р аб И ф
в О р о г		в О р о х	в О р о	в О р а г	в О р а х



## Інші застосування алгоритму

---

Перетворення чисел, символів і скорочень на орфографічний текст (приклад правил з фару)



## Основна ідея лінгвістичної моделі (нагадування)

Імовірність речення із  $l$  слів  $W_{0,l} = (w_1, w_2, \dots, w_l)$  розкладається на добуток умовних імовірностей:

$$P(w_1, w_2, \dots, w_l) = \prod_{i=1}^l P(w_i / w_1, w_2, \dots, w_{i-1}) = \prod_{i=1}^l P(w_i / W_{0,i-1}),$$



## Основна ідея лінгвістичної моделі

Обмежуємо контекст:

$$P(w_1, w_2, \dots, w_l) \cong \prod_{i=1}^l P(w_i / w_{i-n+1}, \dots, w_{i-1}) = \prod_{i=1}^l P(w_i / W_{i-n, i-1}),$$

де  $n$  – порядок моделі.

При  $n = 1$  модель вироджується до відсутності елементарних контекстів, тобто розглядаються **уні-грами**.

При  $n = 2$  отримуємо **бі-грамну** модель, коли в елементарний контекст входить граматична функція одного попереднього слова.

Збільшуючи далі параметр  $n$ , ми розширюємо елементарний контекст, підвищуючи точність моделі.



## Основна ідея лінгвістичної моделі

Враховуючи:

$$\hat{P}(W_{0,l}) = \prod_{i=1}^l P(w_i / W_{i-n,i-1}),$$

Оцінюємо ймовірність послідовності слів при надходженні наступного слова  $w_{i+1}$

$$\hat{P}(W_{0,i+1}) = \hat{P}(W_{0,i})P(w_{i+1} / W_{i-n,i})$$



## Основна ідея лінгвістичної моделі

Оцінка параметрів лінгвістичної моделі

$$P(w_i/W_{i-n,i-1}) = \begin{cases} \alpha(W_{i-n,i-1}) & : & c(W_{i-n,i}) = 0 \\ d_{c(W_{i-n,i})} \frac{c(W_{i-n,i})}{c(W_{i-n,i-1})} & : & 1 \leq c(W_{i-n,i}) \leq k \\ \frac{c(W_{i-n,i})}{c(W_{i-n,i-1})} & : & c(W_{i-n,i}) > k \end{cases}$$



## Основна ідея лінгвістичної моделі

---

Розбиття слів на класи  $g = G(w)$





## Узгодження розшифрованих елементів тексту

---

Класи, яким належить слово  $w$ :

$$G(w) = \{g^1(w), g^2(w), \dots, g_n(w)\}$$

Задаються на основі граматичних типів (частина мови, число, відмінок, рід, час, особа ...)



## Узгодження розшифрованих елементів тексту

Зафіксуємо по класу з множини класів  $G(w)$  для кожного слова  $w$ :

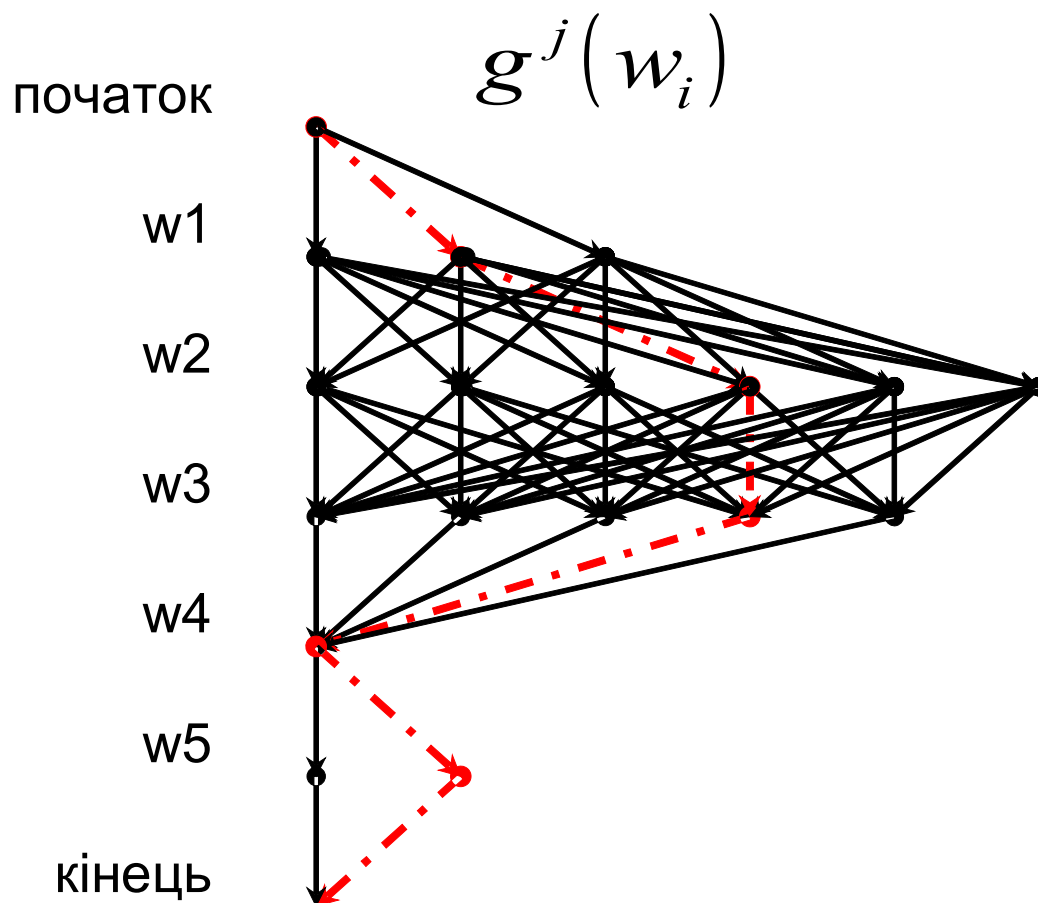
$$G_{0,l}^j = (g^j(w_1), g^j(w_2), \dots, g^j(w_l)).$$

Оцінюємо ймовірність послідовності класів слів при надходженні наступного слова  $w_{i+1}$  за відомої  $\hat{P}(G_{0,i}^k)$  для  $j$ -го класу слова  $w_{i+1}$ :

$$\hat{P}(G_{0,i}^j) = \hat{P}(G_{0,i}^k) P(g^j(w_{i+1}) / G_{i-n,i}^k)$$



## Узгодження розшифрованих елементів тексту



Літня школа-семінар “Усномовні технології та проблеми створення корпусів” Житомир’2010



## Поділ на синтагми (інтонаційні групи)

---

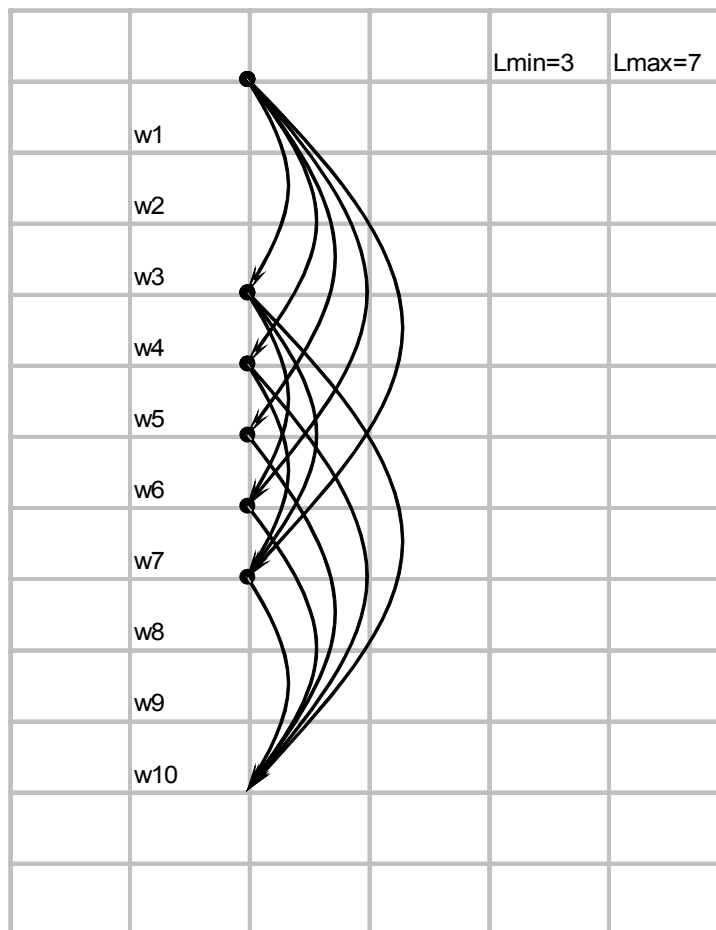
$IG(w(i))$  повертає номер синтагми слова  $w(i)$  у реченні  $W$

$$P(IG(w(i)) = j \mid IG(w(i-1)) = j-1)$$

$L_{min}$ ,  $L_{max}$  – найкоротша і найдовша тривалість синтагми



## Поділ на синтагми (інтонаційні групи)



Літня школа-семінар “Усномовні технології та проблеми створення корпусів” Жуків’2010



## Поділ на синтагми

---

$g(i)$  – граматична форма слова  $w(i)$

Тоді

$$P(IG(w(i)) = j \mid IG(w(i-1)) = j-1)$$

заміняємо на

$$P(IG(g(i)) \neq IG(g(i-1)))$$



Спасибі за увагу!

---

[robeiko@uasoiro.org.ua](mailto:robeiko@uasoiro.org.ua)

[mykola@uasoiro.org.ua](mailto:mykola@uasoiro.org.ua)



Літня школа-семінар “Усномовні технології та проблеми створення корпусів” Жуків’2010