

УДК 004.934.1'1

**В.Ю. Шелепов, А.В. Ниценко**

Институт информатики и искусственного интеллекта ДонНТУ(Украина),  
 Институт проблем искусственного интеллекта НАН и МОН Украины  
 Украина, 83048, м. Донецк, ул. Артема, 118-б

## К проблеме распознавания слитной речи

**V.Ju. Shelepov, A.V. Nicenko**

*Institute of Informatics and Artificial Intelligence DonNTU (Ukraine),  
 Institute of Artificial Intelligence MES of Ukraine and MAS of Ukraine, c. Donetsk  
 Ukraine, 83048, c. Donetsk, Artema st., 118-b*

## *On the Problem of Continuous Speech Recognition*

**В.Ю. Шелепов, А.В. Ниценко**

Інститут інформатики й штучного інтелекту ДонНТУ (Україна),  
 Інститут проблем штучного інтелекту МОН України і НАН України, м. Донецьк  
 Україна, 83048, м. Донецьк, вул. Артема, 118-б

## До проблеми розпізнавання зливої мови

В статье описывается предлагаемый авторами метод распознавания слитно произносимых фраз, состоящих из слов заданного словаря. Метод основан на использовании разработанного авторами механизма сегментации речевого сигнала и алгоритме нахождения первого слова, который авторы обозначают как «принцип минимума DTW-расстояния». Используется ранее предложенная авторами методика построения эталонов слов из дифонов и некоторая модификация метода DTW, дающая заметный выигрыш в скорости и объеме необходимой памяти. Практическая реализация метода требует также уточнения алгоритмов определения границ записанного речевого отрезка и расширения множества правил автоматического транскрибирования. Этому посвящены два первых раздела статьи.

**Ключевые слова:** начало и конец речи, транскриптор, слитная речь, сегментация, дифон, алгоритм DTW, выделение первого слова, принцип минимума DTW-расстояния.

Author's recognition of continuous speech method is describing in the article. Phrase is consisting of words of given vocabulary. The method is based on using author's mechanism of segmentation speech signal and algorithm of searching the first word, which we name *minimal DTW-distant principle*. We use earlier proposed by authors procedure of construction patterns with the help of diphone base and some modification of method DTW which gives win in speed and memory. Practical realization of method demands of more precise determination of speech boundaries and enlargement quantity rules of automatic transcription system. This is the subject of two first sections of the article.

**Key words:** beginning and end of speech, transcription system, continuous speech, segmentation, diphone, DTW-algorithm, the first word separation, minimal DTW-distant principle.

У статті описується запропонований авторами метод розпізнавання зливо вимовлених фраз, що складаються зі слів заданого словника. Метод заснований на використанні розробленого авторами механізму сегментації мовного сигналу й алгоритму знаходження першого слова, що автори позначають як «принцип мінімуму DTW-відстані». Використовується раніше запропонована авторами методика побудови еталонів слів з дифонів і деяка модифікація методу DTW, що дає помітний вииграш у швидкості й обсязі необхідної пам'яті. Практична реалізація методу вимагає також уточнення алгоритмів визначення границь записаного мовного відрізка й розширення множини правил автоматичного транскрибування. Цьому присвячені два перших розділи статті.

**Ключові слова:** початок і кінець мови, транскриптор, злита мова, сегментація, дифон, алгоритм DTW, виділення першого слова, принцип мінімуму DTW-відстані.

## 1 Видоизмененный алгоритм определения начала и конца речевого отрезка

Описываемый ниже алгоритм продолжает тему работы [1] и ориентирован на снижение влияния шума микрофона и звуковой карты.

Используется 8-битная запись с частотой 22050 Гц. По нажатии кнопки записи записываются последовательные отрезки звука по 300 отсчетов (окна). Для каждого из них вычисляется отношение  $V/C$ , где  $V = \sum_{i=0}^{298} |x_{i+1} - x_i|$  – численный аналог полной вариации,  $C$  – количество точек постоянства, то есть таких моментов времени, что в следующий момент величина сигнала остается той же самой. Берется среднее этого отношения по первым 10 окнам. Назовем эту величину «текущий StartPorog». Она характеризует верхний порог «молчания». Ждем момента, когда этот порог будет превышен не менее 5 раз подряд. Возвращаемся на 20 окон назад (начальный запас) и, начиная с этого момента, заносим записываемые отсчеты в буфер 1. Тем самым начинается запись того, что мы предполагаем речью. Определим «текущий EndPorog» как пятикратный текущий StartPorog. Заполнение буфера 1 продолжается до момента, после которого величины  $V/C$  на протяжении 10 тысяч отсчетов будут меньше, чем текущий EndPorog. В него заносятся также упомянутые 10 тысяч отсчетов (запас в конце). Таким образом, запись предполагаемого речевого отрезка останавливается. Отметим, что при каждой записи вычисляются новые значения величин «текущий Start Porog» и «текущий EndPorog».

Записанное проверяется на наличие речи с использованием квазипериодичности ([2]). Если наличие речи обнаруживается, содержимое буфера 1 передается в буфер 2.

Записанный речевой отрезок сегментируется ([3]). Ввиду сказанного выше, сегментация будет начинаться и заканчиваться отрезком паузы (маркировка символом  $P$ ). Наличие этого отрезка в конце позволяет определять, предшествует ли ему гласный ( $W$ ) или звонкий согласный ( $C$ ). Если заключительному  $P$ -отрезку непосредственно предшествует шипящий звук ( $F$ ), алгоритм сегментации также позволяет его обнаружить.

Шум звуковой карты и микрофона может исказить информацию о границах речи. В связи с этим производится уточнение левой границы речевого отрезка. Для этого все записанное подвергается 100-кратному сглаживанию. При этом начальный отрезок молчания превращается в функцию времени, близкую к постоянной (значение этой постоянной определяется величиной первого отсчета, записанного в буфер 2). Считаем, что речь начинается с момента, когда отклонение от этой постоянной превышает порог  $p_1$  (у нас это 10). Отмечаем этот момент в сигнале с помощью метки. Столь сильное сглаживание может «обрезать» начальный шипящий или часть звонкого согласного. Поэтому, если сегментация, произведенная выше, обнаруживает в начале отрезок шипящего или звонкого согласного, метка начала речи при необходимости сдвигается влево, в положение начала шипящего или согласного. Символ  $P$  в начале записи убирается.

Аналогичным образом производится уточнение правой границы записанного речевого отрезка.

Определяется наличие или отсутствие в конце речи глухого взрывного звука ( $L$ ,  $K$ ,  $T$  или их мягкие варианты). Для этого подсчитывается расстояние (количество отсчетов) между последней меткой  $P$  и уточненной меткой конца сигнала. Если оно превышает некоторый порог  $p_2$  (у нас это 2500), то считаем что в конце речи есть глухой взрывной и оставляем заключительный отрезок с маркировкой  $P$  у его левой границы. Если это расстояние меньше  $p_2$ , то заключительный  $P$ -отрезок убирается вместе с маркировкой метки  $P$  и эта метка считается истинным концом сигнала.

## Пример 1

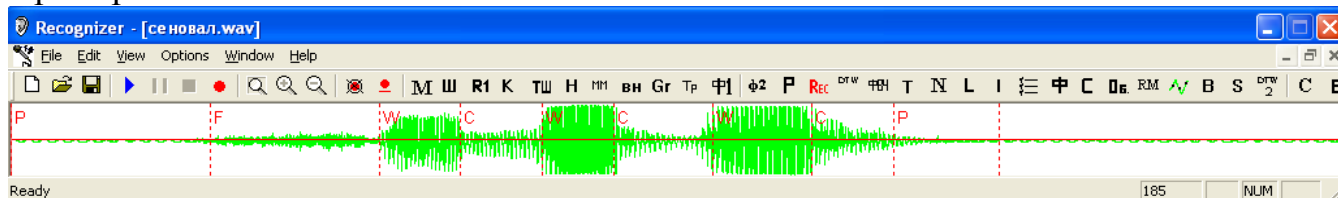


Рисунок 1 – Результат предварительной записи слова «Сеновал» с сегментацией

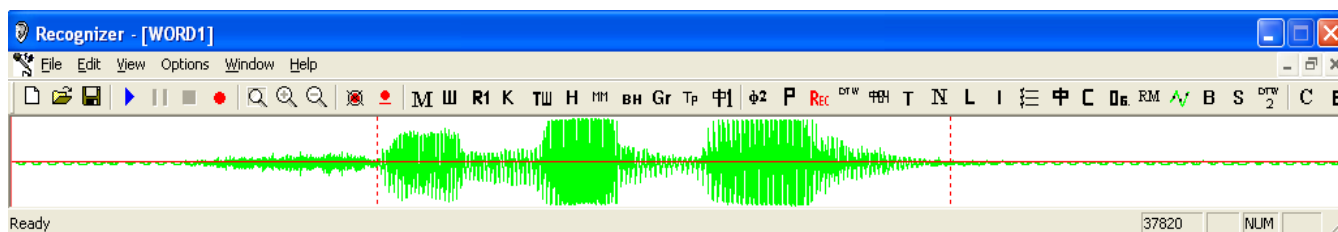


Рисунок 2 – Границы в том же слове после сглаживания

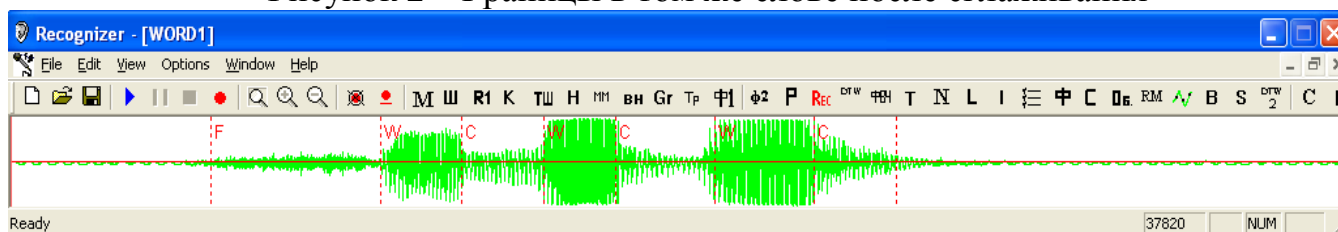


Рисунок 3 – Окончательный результат записи и сегментации слова «Сеновал»

## Пример 2

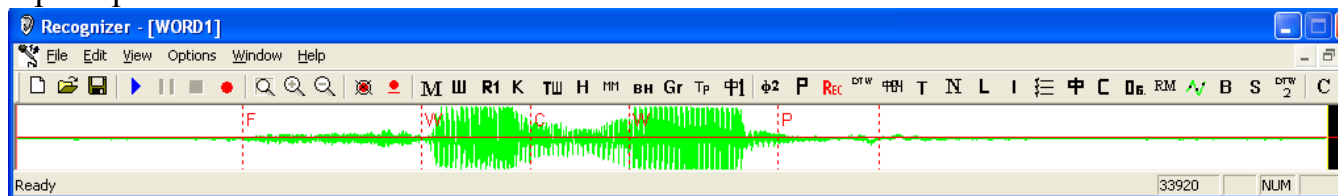


Рисунок 4 – Окончательный результат записи и сегментации слова «Салат»

Отметим в заключение, что использование упомянутой выше проверки на наличие речи позволяет организовать механизм автоматической записи, когда программа, записав слитный речевой отрезок, ожидает и записывает следующий, и пользователю нет необходимости каждый раз нажимать кнопку записи.

## 2 Расширенный транскриптор

Вновь появившиеся в связи с использованием дифонов возможности для различения звонких взрывных между собой, такие же возможности относительно глухих взрывных, твердых и мягких звуков, а также возникновение новых фонетических ситуаций на стыке слов при распознавании слитной речи, потребовало существенного расширения множества правил, заложенных в нашем автоматическом транскрипторе.

В качестве транскрипционных знаков для гласных звуков использованы в основном соответствующие русские буквы. Исключение составляют символы *w*, *q* для ударных *Е*, *Я* соответственно. Об особенностях этой ситуации сказано ниже. Твердые русские согласные транскрибируются также русскими буквами, а соответствующие мягкие согласные – аналогичными латинскими буквами. Исключения: значком @ обозначается мягкое *П*, значком \$ – мягкое *Ж*, значком & – южнорусское (украинское) *Г*, значком + обозначается слитный звук  $\overline{д' ж'}$  (звонкая параллель *Ч*), значком % – слитный звук  $\overline{дз}$  (звонкая параллель *Ц*).

Несколько предварительных слов о принципах, которых мы придерживаемся. Прежде всего, специалистам по распознаванию образов известно, что, как правило, увеличение числа классов распознавания ведет к снижению надежности распознавателя. Поэтому мы сознательно отказываемся от попыток тонкого распознавания аллофонов гласных фонем и используем для них единые транскрипционные символы *A, И, O, У, Э*, за одним исключением. В русском письме в том случае, когда за мягким согласным следует гласный, мягкость отражается путем выбора буквы для гласного: например, *A* после твердого *Д* в слове «да» и *Я* после мягкого *Д* в слове «дядя». Фонетически эти случаи отличаются очень сильно. Поэтому мы сохраняем в качестве транскрипционных знаков буквы *E, Ě, Ю, Я*.

Транскриптор реализован как программа, заменяющая одни символы другими в соответствии с правилами, содержащимися в управляющем файле.

Вот его содержание на сегодняшний день:

1) ~ ~

2) ~e=jе, ~ě=jě, ~ю=јю, ~я=јя, ~\e=j\е, ~\ě=j\ě, ~\ю=ј\ю, ~\я=ј\я

3) ~ =

4) " =

. =

, =

? =

! =

5) ого#=ова, его#=ева,

б) \асч=\аш, \исч=\иш, \осчи=\оши, в\есч=в\еш, досч\и=дош\и, исч\ез=иш\ез, исчез=ищез, насч=наш, обсч=общ, пересч\и=переш\и, пересчи=перещи, песч=пеш, пр\осч=пр\ош, расч\ёс=расч\ёс=раш\ёс, сч\ас=щ\ас, счас=щас, сч\ат=щ\ат, счето=щето, сч\ёт=щ\ёт, сч\ит=сч\ит=щ\ит, тсчит=тщит, счит=счит=щит, тысяч=тыщ,

7) легк=лехк, лёгк=лĕхк, м\ягк=мяхк,

8) здн=зн, дц=ц, тц=щц, тч=чч, жч=щ, зж=жж, сж=жж, сш=шш, стс=сс, стьс=сц, ндс=нс, нтс=нс, стн=сн, стц=сц, здц=сц,

9) \o=1, o=a, l=\o, \e=w, \я=q,

10) ъе=jе, ъě=jě, ъю=јю, ъя=јя, ъа=јя, ъе=jе, ъě=jě, бэ=jе, бю=јю, бя=јя, ъw=јw, ъ\ě=ј\ě, ъ\ю=ј\ю, ъq=јq, ъw=јw, ъ\ě=ј\ě, ъ\o=ј\ě, ъ\ю=ј\ю, ъq=јq, йа=йя, йу=йю, \а=йq,

11) ае=aje, аю=ajю, ая=ajя, ее=eje, ею=eјю, ея=eјя, ёе=ёје, ёю=ёјю, ёя=ёјя, ие=ije, ию=ijю, ия=ijя, ое=oje, ою=oјю, оя=oјя, уе=uje, ую=ujю, уя=ujя, ые=ыје, ью=ыјю, ья=ыјя, эе=эје, эю=эјю, эя=эјя, юе=юје, юю=юјю, юя=юјя, яе=яје, яю=яјю, яя=яјя, we=wje, wю=wјю, wя=wјя, qe=qje, qю=qјю, qя=qјя, aw=ajw, a\ě=aj\ě, a\ю=aj\ю, aq=ajq, ew=eјw, e\ě=eј\ě, e\ю=eј\ю, eq=eјq, ěw=ěјw, ě\ě=ěј\ě, ě\ю=ěј\ю, ěq=ěјq, iw=ijw, i\ě=ij\ě, i\ю=ij\ю, iq=ijq, ow=ajw, o\ě=aj\ě, o\ю=aj\ю, oq=ajq, uw=ujw, u\ě=uj\ě, u\ю=uj\ю, uq=ujq, ыw=ыјw, ы\ě=ыј\ě, ы\ю=ыј\ю, ыq=ыјq, эw=эјw, э\ě=эј\ě, э\ю=эј\ю, эq=эјq, юw=юјw, ю\ě=юј\ě, ю\ю=юј\ю, юq=юјq, яw=яјw, я\ě=яј\ě, я\ю=яј\ю, яq=яјq, ww=wјw, w\ě=wј\ě, w\ю=wј\ю, wq=wјq, qw=qјw, q\ě=qј\ě, q\ю=qј\ю, qq=qјq,

12) #e=jе, #ě=jě, #ю=јю, #я=јя, #w=јw, #\ě=ј\ě, #\ю=ј\ю, #q=јq,

~e=~je, ~ě=~jě, ~ю=~јю, ~я=~јя, ~w=~јw, ~\ě=~ј\ě, ~\ю=~ј\ю, ~q=~јq, стьд=zd,

13) бѣ=b, бј=bј, бе=be, бě=bě, би=би, бю=бю, бя=бя, бq=bq, бw=bw, б\ě=b\ě, б\и=b\и, б\ю=b\ю, бb=bb, вѣ=v, вј=vј, ве=ve, ви=ви, вю=вю, вя=вя, вq=vq,

vw=vw, vë=vë, v\ë=v\ë, v\и=v\и, v\ю=v\ю, vv=vv, гь=g, гj=gj, ге=ge, ги=ги, гю=гю, гя=гя, гq=гq, gw=gw, гë=gë, г\ë=g\ë, г\и=g\и, г\ю=g\ю, гg=gg, дь=d, dj=dj, де=de, ди=ди, дю=дю, дя=дя, dq=dq, dw=dw, дë=dë, д\ë=d\ë, д\и=d\и, д\ю=d\ю, dd=dd, зь=z, зj=zj, зе=ze, зи=зи, зю=зю, зя=зя, зq=zq, zw=zw, зë=zë, з\ë=z\ë, з\и=z\и, з\ю=z\ю, zz=zz, кь=k, kj=kj, ке=ke, кë=kë, ки=ки, кю=кю, кя=кя, kq=kq, kw=kw, к\ë=k\ë, к\и=k\и, к\ю=k\ю, kk=kk, ль=l, lj=lj, ле=le, ли=ли, лю=лю, ля=ля, лq=lq, лw=lw, л\ë=l\ë, лë=lë, л\и=l\и, л\ю=l\ю, ll=ll, мь=m, mj=mj, ме=me, ми=ми, мю=мю, мя=мя, mq=mq, mw=mw, мë=m\ë, м\ë=m\ë, м\и=m\и, м\ю=m\ю, mm=mm, нь=n, nj=nj, не=ne, ни=ни, ню=ню, ня=ня, nq=nq, nw=nw, нë=në, н\ë=n\ë, н\и=n\и, н\ю=n\ю, nd=nd, пь=@, pj=@j, пе=@e, пë=@ë, пи=@и, пю=@ю, пя=@я, pq=@q, pw=@w, п\ë=@\ë, п\и=@\и, п\ю=@\ю, п@=@@, рь=r, рj=rj, ре=re, рë=rë, ри=ри, рю=рю, ря=ря, rq=rq, rw=rw, рë=rë, р\ë=r\ë, р\и=r\и, р\ю=r\ю, rr=rr, сь=s, sj=sj, се=se, си=си, сю=сю, ся=ся, cq=sq, cw=sw, сë=së, с\ë=s\ë, с\и=s\и, с\ю=s\ю, ct=st, фь=f, fj=fj, фе=fe, фи=фи, фю=фю, фя=фя, fq=fq, fw=fw, фë=fë, ф\ë=f\ë, ф\и=f\и, ф\ю=f\ю, ff=ff, xj=hj, xe=he, хи=хи, хю=хю, хя=хя, xq=hq, xw=hw, хë=hë, х\ë=h\ë, х\и=h\и, х\ю=h\ю, хh=hh, ть=t, tj=tj, те=te, тë=të, ти=ти, тю=тю, тя=тя, tq=tq, tw=tw, тë=të, т\ë=t\ë, т\и=t\и, т\ю=t\ю, tt=tt, zd=zd, zl=zl,

14) лzn=л2n, рzn=р2n, zn=zn, 2=з, ннщ=нщ, nn=nn, nt=nt, нч=нч, нщ=нщ, ccl=ssl, cl=sl, ccn=ssn, cn=sn, лct=л3t, pct=p3t, ct=st, 3=c, cs=ss,

15) ь=

16) б#=п, в#=ф, г#=к, д#=т, ж#=ш, з#=с, б#=@, в#=f, д#=t, з#=s,

17) бк=пк, бп=пп, бс=пс, бт=пт, бф=пф, бх=пх, бц=пц, бш=пш, вк=фк, вп=фп, вс=фс, vt=ft, вф=фф, vx=fx, вц=фц, вш=фш, гк=кк, гп=кп, гс=кс, гт=кт, гф=кф, гх=кх, гц=кц, гш=кш, dk=tk, дп=тп, dc=tc, dt=tt, дф=тф, dx=tx, дц=тц, дш=тш, жк=шк, жп=шп, жс=шс, жт=шт, жф=шф, жх=шх, жц=шц, жш=шш, зк=ск, зп=сп, зс=сс, зт=ст, зф=сф, zx=cx, зц=сц, зш=шш, бк=пк, б@=@@, бs=ps, бt=pt, бf=pf, бh=ph, бч=пч, бщ=пщ, vk=fk, в@=ф@, vs=fs, vt=ft, vf=ff, vh=fh, вч=фч, вщ=фщ, гk=kk, г@=к@, gs=ks, gt=kt, gf=kf, gh=kh, гч=кч, гщ=кщ, dk=tk, зtk=ctk, д@=т@, ds=ts, dt=tt, df=tf, dh=th, wздч=wщ, здч=щч, дч=тч, дщ=тщ, жk=шк, ж@=ш@, жs=шs, жt=шт, жf=шf, жh=шh, жч=шч, жщ=шщ, zk=ck, з@=c@, zs=ss, зt=ct, зf=cf, zh=ch, зч=сч, зщ=щщ, bk=@k, bp=@п, bc=@c, bt=@т, bf=@ф, bx=@x, бц=@ц, бш=@ш, vk=fk, vp=fp, vc=fc, vt=ft, vф=фф, vx=fx, вц=фц, vш=фш, dk=tk, dp=tp, dc=tc, dt=tt, дф=тф, dx=tx, дц=тц, дш=тш, zk=sk, zp=sp, zc=ss, зt=ст, зф=сф, zx=sx, зц=сц, зш=сш, bk=@k, б@=@@, bs=@s, bt=@t, bf=@f, bh=@h, бч=@ч, бщ=@щ, vk=fk, в@=ф@, vs=fs, vt=ft, vf=ff, vh=fh, вч=fч, вщ=fщ, dk=tk, д@=т@, ds=ts, dt=tt, df=tf, dh=th, дч=тч, дщ=тщ, zk=sk, з@=s@, zs=ss, зt=st, zh=sh, зч=sч, зщ=щщ,

18) kb=gb, kg=gg, kd=gd, kz=gz, pb=bb, pg=bg, pd=bd, pz=bz, cb=zb, cg=zg, cd=zd, cz=zz, fb=bb, fg=bg, fd=bd, fz=bz, xb=&b, xg=&g, xd=&d, xz=&z, цb=%b, цг=%г, цд=%д, цж=%ж, цз=%з, чb=+б, чг=+г, чд=+д, чж=+ж, чз=+з, шb=жб, шг=жг, шд=жд, шж=жж, шз=жз, щb=\$б, щг=\$г, щд=\$д, щж=\$ж, щз=\$з, kb=gb, kg=gg, kd=gd, kz=gz, pb=bb, pg=bg, pd=bd, pz=bz, cb=zb, cg=zg, cd=zd, cz=zz, fb=bb, fg=bg, fd=bd, fz=bz, xb=&b, xg=&g, xd=&d, xz=&z, цb=%b, цг=%г, цд=%д, цж=жб, шг=жг, шд=жд, шз=жз, @б=бб, @г=бг, @д=бд, @ж=бж, @з=бз, sb=zb, st=zg, sd=zd, sj=жж, s3=z3, тб=dб, тр=др, тд=дд, тж=dж, тз=dз,

fb=vb, fg=vg, fd=vd, fж=vж, fz=vz, @b=bb, @g=bg, @d=bd, @z=bz, sb=zb, sg=zg, sd=zd, sz=zz, tb=db, tg=dg, td=dd, tz=dz, fb=vb, fg=vg, fd=vd, fz=vz, чb=+b, чg=+g, чd=+d, чz=+z, щb=\$b, щg=\$g, щd=\$d, щz=\$z,

19) же=жэ, жи=жы, жю=жу, жя=жа, жw=ж\э, жё=ж\о, ж\ё=ж\о, ж\и=ж\ы, ж\ю=ж\у, жq=ж\а, ше=шэ, ши=шы, шю=шу, шя=ша, шw=ш\э, ш\ё=ш\о, ш\и=ш\ы, ш\ю=ш\у, шq=ш\а, це=цэ, ци=цы, цю=цу, ця=ца, цw=ц\э, ц\ё=ц\о, ц\и=ц\ы, ц\ю=ц\у, cq=ц\а, ча=чя, чу=чю, чэ=че, ч\а=чq, ч\о=чё, ч\у=чю, ч\э=чw, ща=щя, шу=шю, шэ=ще, щ\а=щq, щ\о=щё, щ\у=щю, щ\э=щw,

20) лл#=#л, мм#=#м, нн#=#н,

21) \=

Поясним приведенный перечень правил. Каждое из них записано в виде двух или более частей, соединенных знаком =. Если упомянутых частей две, слева стоят исходные символы буквенной записи слова, справа – символы которыми они заменяются в транскрипции. Значок \ означает ударение. Машина, транскрибируя слово, последовательно ищет вхождение левой части очередного правила, и если таковое обнаруживается, заменяет его правой частью. Если упомянутых частей больше двух, создается соответствующее число вариантов транскрипции: вариант, соответствующий второй части, вариант, соответствующий третьей части равенства, и т.д.

Для удобства читателя в данном тексте правила разбиты на группы, которые пронумерованы. Рекомендуется внести в управляющий файл эти группы в порядке номеров, не меняя порядка правил в группах, поскольку порядок замен, очевидно, важен.

Правило первой группы введено исключительно для наглядности. Оно временно заменяет пробел знаком ~.

Вторая группа описывает произношение *Е, Ё, Ю, Я* после пробела (начало слова).

Правило третьей группы убирает значок ~, а вместе с ним из транскрипции слитной речи уходят пробелы между словами.

Четвертая группа удаляет из транскрипции знаки препинания.

Пятая группа описывает произношение окончаний в родительном падеже прилагательных типа «нового», «синего».

Шестая группа служит для транскрибирования сочетания «СЧ» в различных ситуациях. Сочетание СЧ, которое в слове «считать» от слова «счёт» звучит, как Щ, даёт СЧ в омониме, обозначающем чтение с какого-то носителя; это порождает два варианта транскрипции.

Седьмая группа предназначена для описания произношения в словах типа «легко».

Восьмая – отражает произносительную норму в словах типа «мужчина» и некоторые фонетические правила, связанные с неизносимыми согласными.

Девятая группа служит для транскрибирования гласных *О, Е, Я*, когда они стоят в ударной позиции. Поскольку транскриптор работает по принципу замены, приходится предварительно переименовывать ударное *О*, а затем возвращать ему прежнее обозначение. Введение специальных обозначений для ударных *Е, Я* связано с тем, что только они имеют совершенно определенное произношение. В безударном варианте они произносятся различными носителями языка по-разному. Для так называемой «младшей нормы» (более молодое поколение москвичей) они ближе к *И*, у сибиряков и в сценической речи – ближе к *Е, Я*.

Десятая группа правил отражает фонетическую роль мягкого и твердого знаков перед *Е, Ё, Ю, Я*. Их наличие приводит при произношении к появлению согласного *ј*.

Одиннадцатая группа отражает произношение сочетаний гласных с гласными *Е, Ё, Ю, Я*.

Двенадцатая группа описывает произношение  $E, \dot{E}, Ю, Я$ , если с них начинается произносимый слитный речевой отрезок (# – знак начала и конца; в транскрибируемом тексте его проставлять не надо).

Тринадцатая и четырнадцатая группы связаны с обозначением в русском письме мягкости согласных и нейтрализацией твердых и мягких фонем.

Пятнадцатая – удаление мягкого знака, который уже сыграл свою роль.

Шестнадцатая группа – оглушение звонкой согласной в конце произносимого слитного речевого отрезка.

Семнадцатая группа – оглушение звонкой согласной перед глухой взрывной, шипящей и аффрикатами  $Ц, Ч$ .

Восемнадцатая – озвончение глухих согласных перед звонкими согласными.

Девятнадцатая группа отражает влияние твердого и мягкого согласного на последующий гласный.

Двадцатая – особенность произнесения удвоенных согласных в конце слова.

Отметим, что мы опробовали ряд достаточно успешных методов автоматического определения ударения в слове. Решение до конца этой трудной проблемы, очевидно, сильно сократило бы число слов – кандидатов на распознавание. Наш транскриптор пока убирает знак ударения (группа 21), но, в расчете на его автоматическое определение в звучащем слове, делает это лишь в самом конце.

Отметим, что мы включили в приведенный перечень лишь те правила, которые обусловлены русской фонетикой, и оставили за его пределами некоторые правила, порождаемые особенностями нашей сегментации. Например, сонорные согласные на конце слова после глухих взрывных сегментируются как гласные звуки (идентификатор  $W$ ). Эти дополнительные правила включаются в отдельный транскриптор, используемый нами для создания файлов слов с транскрипцией широкой фонетической классификации ( $W, C, F, P$ ), которые в данной работе не используются. Отметим также, что ряд вышеприведенных сочетаний не встречается в отдельных словах, но встречается в слитной речи.

Наконец, прежде, чем транскрибировать по указанным правилам, компьютер обращается к файлу исключений, в котором описываются процедуры транскрибирования целых слов, например,

чт\o=што, ог\o=ого.

### 3 Алгоритм с DTW-эталонами, создаваемыми из дифонов. Использование таблицы расстояний.

#### Дерево эталонов. Размеры DTW-матрицы

Мы применяем для распознавания ставший уже классическим алгоритм Т.К. Винцока, известный под названием алгоритма DTW (его описание в [1], также [2]). При этом мы используем свои вектора признаков, связанных с относительными частотами длин полных колебаний на речевых отрезках в 368 отсчетов ([4]). Эталоны слов распознаваемого словаря формируются из эталонов дифонов, полная база которых в объеме приблизительно полуторных тысяч создается для каждого диктора заранее (2 – 3 часа работы, [4]). Отметим, что создание такой базы в дальнейшем избавляет пользователя от необходимости создавать какие-либо эталоны голосом.

Под дифоном, соответствующим межфонемному переходу внутри слова, будем понимать участок стандартной длины: 3 окна в 368 отсчетов слева от метки между звуками и 3 таких же окна справа от той же метки. Эталон дифона – набор 6-и соответствующих векторов. Кроме того, мы используем участок в 3 окна в начале слова и участок в 3 окна – в конце слова, условно называя их соответственно начальным и конечным полудифоном слова (переход от молчания к речи и наоборот). Все вектора, входящие

в эталоны дифонов, играют роль кодовых векторов и образуют кодовую книгу  $B$ . Все эталоны дифонов нумеруются, нумеруются также все кодовые вектора.

Каждое слово словаря автоматически транскрибируется, по транскрипции строится цепочка имен дифонов. Каждое из них заменяется эталоном соответствующего дифона. Полученная цепочка векторов образует эталон слова ([4]).

На самом деле мы не создаем и не храним перечень эталонов слов словаря в виде статического списка. Словарь эталонов слов реализуется в виде дерева дифонов, использование которого существенно ускоряет процесс распознавания. Дерево создается при первоначальной загрузке текстового словаря. Дифоны представлены в дереве своими номерами. Эталон каждого слова представляется в виде ветви этого дерева. Если несколько ветвей имеют общую часть, то вычисления, заполняющие соответствующую часть DTW-матрицы, выполняются только один раз.

Уровни дерева соответствуют позициям дифонов в слове. Каждый узел в рамках каждого уровня представляет собой номер дифона, находящегося в слове на соответствующей позиции. Вершины, соответствующие конечным дифонам слов, помечаются как концы соответствующих слов (в узле записывается порядковый номер соответствующего слова в словаре). Если узел не конечный, то записывается значение -1. Максимальная глубина дерева соответствует максимальной длине (выраженной в количестве дифонов) слова в словаре.

Процесс распознавания строится следующим образом. Распознаваемое слово автоматически сегментируется и затем подвергается так называемой межфонемной обработке: удаляются стационарные части составляющих звуков и оставляются лишь дифоны в окрестностях межзвуковых меток (межфонемные переходы). Затем создается представление слова в виде набора  $N$  векторов признаков и строится таблица  $D$  расстояний этих векторов до всех векторов кодовой книги  $B$ . Далее вычисляются DTW-расстояния рассматриваемого слова до всех эталонов слов путем рекурсивного обхода дерева эталонов «в глубину». Вначале просматриваем корень дерева, а затем спускаемся по ветви, пока не достигнем вершины, помеченной как конец слова. После того, как достигнут конец слова, возвращаемся назад вдоль пройденного пути, пока не найдем вершину, у которой есть еще не посещенный сосед, а затем двигаемся в новом обнаруженном направлении. Процесс оказывается завершенным, когда мы вернулись в корень дерева, а все примыкающие к нему вершины уже оказались посещенными.

При прохождении ветвей дерева, по номерам дифонов строится цепочка соответствующих им номеров векторов, образующих эталон слова. Двигаясь в глубину, добавляем в цепочку номера, соответствующие пройденным узлам, а при движении назад они удаляются из нее. Достигнув узла, являющегося концом очередного слова, вычисляем DTW-расстояние от построенной цепочки векторов (эталона данного слова) до цепочки векторов распознаваемого сигнала. При этом расстояния между векторами берутся из таблицы  $D$ . В процессе вычисления расстояний матрица DTW не переписывается полностью, а обновляются только столбцы, соответствующие новым кодовым векторам, номера которых добавлены в цепочку после возврата назад по окончании предыдущего этапа.

Таким образом, достигается очень значительный выигрыш как в скорости распознавания, так и в объеме необходимой памяти.

В заключение отметим, что мы работаем с квадратными DTW-матрицами переменного размера: если эталон слова содержит  $a$  векторов, а распознаваемое слово содержит  $b$  векторов, то мы строим DTW-матрицу размера  $\sqrt{a^2 + b^2}$  (Г.В. Дорохина).



## 4 Основные принципы предлагаемого подхода к распознаванию слитной речи

Пусть у нас есть несколько слитно произнесенных фраз. Наша программа автоматически затранскрибирует их и создаст для каждой из них эталон из дифонов, игнорируя пробелы между словами. После этого их можно распознавать между собой теми же методами, что и отдельно произносимые слова. Но если рассматривать множество произвольных фраз, то их бесконечно много и, очевидно, следует добиваться их распознавания путем распознавания слов, из которых они состоят. Тогда основная сложность – выделение в речевом сигнале отрезков, отвечающих отдельным словам. Иначе говоря, мы должны научиться определять, где заканчивается одно слово и начинается другое. Предлагаемый ниже метод основан на использовании вышеупомянутой сегментации. Весь рассматриваемый речевой отрезок автоматически разбивается на сегменты, отвечающие отдельным звукам, и границы между словами следует искать среди конечного множества полученных границ между звуками.

Мы начинали с распознавания пар слитно произносимых слов. Распознавая отрезок от начала до первой метки, а затем от первой метки до конца, мы получали пару слов нашего словаря. Затем мы проводили распознавание от начала до второй метки и от второй метки до конца и так далее. Заключительным шагом было распознавание всего речевого отрезка от начала до конца как одного слова. В результате мы получали последовательность гипотетических пар слов (на последнем месте – одно слово). Для каждой из этих пар автоматически строился эталон и результатом распознавания объявлялась пара, до которой DTW-расстояние минимально. Этот алгоритм показал высокую надежность. Но он включал целый набор актов распознавания отдельных гипотетических слов и в результате оказывался довольно долго работающим. Попытка применить аналогичный алгоритм к распознаванию большего числа слитно произнесенных слов ведет к экспоненциальному росту числа распознаваний гипотетических слов, и от нее приходится отказаться.

Тогда мы стали, двигаясь от начала до очередной метки, выводить только последовательность гипотез для первого слова, но с указанием DTW-расстояния до каждой из них. Оказалось, что гипотеза, соответствующая истинному первому слову (и соответствующему истинному отрезку от начала) имеет указанное расстояние, близкое к минимальному. Для слитно произносимых числительных (без фонетических вложений, о которых ниже) результат оказывался точным. Итак, мы приходим к следующему «принципу минимума»: **ПО КРАЙНЕЙ МЕРЕ, ДЛЯ СЛОВАРЕЙ, УДОВЛЕТВОРЯЮЩИХ НЕКОТОРЫМ ОГРАНИЧЕНИЯМ, ПЕРВОЕ СЛОВО ОПРЕДЕЛЯЕТСЯ С ИСПОЛЬЗОВАНИЕМ МЕТОК ИЗ УСЛОВИЯ МИНИМУМА DTW-РАССТОЯНИЯ.**

Понятно, что для распознавания второго слова фразы следует применить описанный метод к части сигнала от конца первого слова до конца речевого отрезка и так далее.

Смысл этого принципа становится понятен, если вспомнить, что алгоритм DTW направлен на минимизацию расстояния сказанного слова до эталона того же слова. Остальные слова в полученном списке на самом деле не звучали и то, что их расстояния до соответствующих эталонов оказались больше, представляется естественным.

Об ограничениях, упомянутых выше. Эксперименты показывают, что к числу таких ограничений нужно отнести следующее. Словарь не должен содержать пар слов,

одно из которых совпадает с началом другого. Точнее, не должно быть таких пар слов, что транскрипция одного из них получается из транскрипции другого приписыванием в конце дополнительных транскрипционных символов. В противном случае, при произнесении более длинного слова такой пары, DTW-расстояние до слова с более короткой транскрипцией может быть меньше.

## Литература

1. Шелепов В.Ю. Новый подход к определению границ речевого сигнала. Проблемы конца сигнала / В.Ю. Шелепов, А.В. Ниценко / Речевые технологии. – Москва, 2012.
2. Шелепов В.Ю. Лекции о распознавании речи / Шелепов В.Ю. – Донецк : ІПШІ Наука і освіта. – 2009. – 192 с.
3. Шелепов В.Ю. Построение системы голосового управления компьютером на примере задачи набора математических формул / В.Ю. Шелепов, А.В. Ниценко, А.В. Жук // Искусственный интеллект. – 2010. – № 3. – С. 259-267.
4. Шелепов В.Ю. О распознавании речи на основе межфонемных переходов / В.Ю. Шелепов, А.В. Ниценко, Г.В. Дорохина // Искусственный интеллект. – 2012. – № 1. – С. 132-139.
5. Винцюк Т.К. Анализ, распознавание и интерпретация речевых сигналов / Винцюк Т.К. – Киев : Наук. Думка, 1987. – 262 с.

## Literatura

1. Shelepov V.Ju. Speech technology / V.Ju. Shelepov, A.V. Nicenko. – Rechevyte tehnologii. – Moskow, 2012.
2. Shelepov V.Ju. Lectures on speech recognition / Shelepov V.Ju. – Donetsk : IAI «Nauka i osvita», 2009. – 192 S.
3. Shelepov V.Ju. Artificial intelegence / V.Ju. Shelepov, A.V. Nicenko, A.V. Zhuk // Iskusstvennyj intellect. – 2010. – №.3. – S. 259-267.
4. Shelepov V.Ju. Artificial intelegence / V.Ju. Shelepov, A.V. Nicenko, G.V. Dorohina // Iskusstvennyj intellect. – 2012. – № 1. – S. 132-139.
5. Vincjuk T.K. Analysis, recognition and interpretation of speech signals / Vincjuk T.K. – Kiev : Naukova dumka, 1987. – 262 S.

### RESUME

*W.Ju. Shelepov, A.V. Nicenko*

### *For the Problem of Continuous Speech Rrecognition*

The first section goes on the subject of article [1]. It describes algorithm determination of speech boundaries more robust relatively noise microphone and sound card. The second section describes automatic transcription system, which contains many new rules of explosive consonants transcription, soft consonants transcription and phonetic situations in boundaries of words.

We apply for recognition some modification of DTW-algorithm, using own system of signs ([4]). We create base of diphones, which contain phone transitions, and build patterns of words pasting together patterns of these diphones. We keep patterns of words in the form of diphones numbers tree. Using this tree we obtain essential saving of speed and memory.

*The main principals of suggested arrangement continuous speech recognition*

1) boundaries of words ought to look for in finite set of marks between sounds (our segmentation builds this set automatically).

2) Let us recognize (as a word of our vocabulary) the part of signal from the beginning to the first mark, then from the beginning to the second mark and so on. Then at least for vocabulary, which satisfy restriction given below, the first word is determined following minimal DTW-distant condition. (minimal DTW-distant principle).

Restriction is the next: vocabulary must not contain so pairs of words, that transcription one of them is the beginning of transcription the other. Otherwise, when we speak more long word, the distant to more shot word may be less.

*Статья поступила в редакцию 06.07.2012.*