

# Информационные технологии и системы

УДК 004.934

В.В. Пилипенко

## ТЕХНОЛОГИЯ РАЗМЕТКИ ЗВУКОВЫХ ФАЙЛОВ С ИСПОЛЬЗОВАНИЕМ НЕТОЧНОГО ТЕКСТОВОГО СОПРОВОЖДЕНИЯ

Описана технология разметки звуковых файлов с использованием неточного текстового сопровождения. Предварительно формируется система распознавания на основе речевых записей, размеченных экспертами. Новые речевые записи распознаются для выяснения временных границ слов. Процедура сравнения ответа распознавания и неточного описания выявляет фрагменты звука, для которых есть точное соответствие. На основе автоматически полученной разметки строится новая, более точная система автоматического многодикторного распознавания спонтанной украинской речи с объемом словаря в 125 тысяч словоформ. Проведенные эксперименты показали пословную точность распознавания в 80 %.

**Введение.** Объемы записанной речевой информации растут вместе с развитием телекоммуникационных технологий. Имеется много приложений, для которых необходима технология, основанная на распознавании произвольной речи. Насущными потенциальными приложениями являются текстовое транскрибирование звуковых записей, поиск аудиоинформации в аудиоархивах, мониторинг медийных источников информации, распознавание тематик речевого файла и т.п.

Современное автоматическое распознавание достигло высоких результатов при распознавании спонтанной многодикторной речи, тем не менее каждая новая область применений требует настройки на новый акустический материал, поскольку одним из основных этапов при построении системы распознавания речи является построение акустической модели для речевого сигнала из данной области применения. Для создания акустических моделей необходимо иметь размеченный звуковой корпус, в котором поставлено точное соответствие между звуком и текстом. Обычно создание такого корпуса требует ручной разметки звукового материала, и соответственно много времени и других ресурсов.

Вместе с этим часто имеются в огромном количестве звуковые записи из телерадиопередач, которые сопровождаются текстами. Эти тексты создаются без учета требования точного соответствия звуку и могут содержать как дополнительные комментарии, не прозвучавшие в звуковом материале, например название передачи или сведения о выступающих, так и пропуски текста, возникшие по разным причинам, например для большей читабельности текста. Такой материал напрямую не годится для обучения системы распознавания, поскольку процедура построения акустических моделей полагается на точное соответствие между звуком и

последовательностями фонем, которые автоматически порождаются по тексту, и в случае когда такого точного соответствия нет, то обученные акустические модели в дальнейшем будут давать значительные ошибки распознавания.

Технология разметки звуковых файлов с использованием неточного текстового сопровождения была разработана и опробована для различных языков [1–4]. Предлагается технология автоматической разметки неточно описанного звукового материала с использованием предварительно созданной системы автоматического распознавания украинской речи.

На предварительном этапе для построения акустических и лингвистических моделей используется аналогичный корпус речи, размеченный экспертами. На следующем шаге созданная предварительная система распознавания речи используется для распознавания нового звукового материала. Процедура сравнения ответа распознавания и текста, который имеется для данного звукового материала, выявляет те фрагменты, для которых есть точное соответствие между звуком и текстом. При наличии большого количества звуковых записей таких фрагментов может оказаться достаточно много для построения акустических моделей, наилучшим образом пригодных для распознавания аналогичных речевых записей.

**Постановка задачи распознавания слов.** Представим произнесенную последовательность слов как последовательность векторов или наблюдений  $O$

$$O = o_1, o_1, \dots, o_T, \quad (1)$$

где  $o_t$  — вектор речи, наблюдаемый в момент времени  $t$ . Проблему распознавания слов можно рассматривать как результат вычисления

$$\arg \max_{\{w\}} \{P(w/O)\}, \quad (2)$$

где  $\{w\}$  — искомая последовательность слов, а поиск максимума производится по всем возможным словам в последовательности. Использование правила Байеса дает способ вычисления вероятности:

$$P(w/O) = \frac{P(O/w)P(w)}{P(O)}. \quad (3)$$

Таким образом, наиболее вероятная произнесенная последовательность слов определяется двумя компонентами  $P(O/w)$  — акустическим и  $P(w)$  — лингвистическим. Вследствие высокой размерности последовательности наблюдений  $O$  прямое оценивание совместной условной вероятности  $P(o_1, o_2, \dots / w)$  не используется на практике. Однако, если сделать предположение о параметрической модели генерирования последовательности слов, такой как марковская модель, оценивание по данным становится возможным, поскольку проблема оценивания условных плотностей  $P(O/w)$  заменяется намного более простой проблемой оценивания параметров марковской модели.

Таким образом, при распознавании речи предполагается, что последовательность наблюдаемых векторов речи, соответствующая

последовательности слов, порождена марковской моделью  $M$ . Марковская модель представляет собой автомат с конечным числом состояний  $X$ , изменяющий свое состояние один раз в каждую единицу времени, и в каждый момент времени  $t$ , когда модель находится в состоянии  $j$ , вектор речи  $o_t$  генерируется исходя из плотности вероятностей  $b_j(o_t)$ . Кроме того, переход от состояния  $i$  к состоянию  $j$  также является вероятностным и совершается под управлением отдельной вероятности  $a_{ij}$ . На практике, однако, известна только последовательность наблюдения  $O$ , а лежащая в основе последовательность состояний  $X$  скрыта. Вот почему такую модель называют *скрытой марковской моделью* (СММ).

Совместная вероятность того, что  $O$  сгенерирована моделью  $M$ , проходящей через последовательность состояний  $X$ , рассчитывается просто как произведение вероятностей перехода  $a_{ij}$  и вероятностей генерирования  $b_j(o_t)$ .

Поскольку последовательность состояний  $X$  неизвестна, требуемое правдоподобие  $P(O/M)$  вычисляется суммированием по всем возможным последовательностям состояний  $X = x(1), x(2), x(3), \dots, x(T)$ , т.е.

$$P(O/M) = \sum_X a_{x(0)x(1)} \prod_{t=1}^T b_{x(t)}(o_t) a_{x(t)x(t+1)}, \quad (4)$$

где  $x(0)$  является начальным состоянием, а  $x(T+1)$  — конечное состояние модели.

Как альтернатива уравнению (4), правдоподобие может быть вычислено приближенно, путем рассмотрения наиболее вероятной последовательности состояний, т.е.

$$P^*(O/M) = \max_X \left\{ a_{x(0)x(1)} \prod_{t=1}^T b_{x(t)}(o_t) a_{x(t)x(t+1)} \right\}. \quad (5)$$

Осуществить прямые вычисления, в соответствии с соотношениями (4) и (5), не очень легко, однако существуют простые рекурсивные процедуры, позволяющие весьма эффективно рассчитать обе величины.

Заметим также, что если уравнение (2) может быть вычислено, тогда проблема распознавания решена. Для заданного множества моделей  $M_i$ , соответствующих последовательности слов  $w_i$ , соотношение (2) решается с использованием (3) и в предположении, что

$$P(O/w_i) = P(O/M_i). \quad (6)$$

При этом, разумеется, предполагается, что параметры  $a_{ij}$  и  $b_j(o_t)$  известны для каждой модели  $M_i$ . Для оценки распределений  $b_j(o_t)$  используются смеси гауссовских распределений, что позволяет небольшим количеством параметров оценить распределения произвольного вида. Для заданного множества примеров обучения, соответствующих конкретной модели, параметры этой модели можно определить автоматически с помощью эффективной рекуррентной процедуры. Таким образом, при условии, что собрано достаточное число представительных образцов каждого слова, может

быть построена СММ, которая неявно моделирует все множество причин изменчивости, свойственной реальной речи.

При распознавании речи для большого словаря использование отдельных моделей для каждого слова становится невозможным, поскольку для этого необходимо, чтобы в обучающей выборке были представлены все слова из словаря распознавания. Обычно модель слова состоит из последовательности моделей фонем, которая задается фонетической транскрипцией слова.

**Процесс оценки параметров акустических моделей.** Чтобы определить параметры акустических моделей, необходимо сделать начальное предположение о том, каковы они могли бы быть. Как только это сделано, можно найти более точные параметры (в смысле максимального правдоподобия) с помощью так называемой рекуррентной процедуры Баума-Уэлша.

Без ограничения общности можно полагать, что компоненты смеси являются специальной формой подсостояния, в котором вероятности перехода являются весами смеси.

Таким образом, важной задачей является оценивание средних и дисперсий СММ, в которой каждое состояние выходных данных представляется единственным гауссовским компонентом:

$$b_j(o_t) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_j|}} e^{-\frac{1}{2}(o_t - \mu_j)' \Sigma_j^{-1} (o_t - \mu_j)}. \quad (7)$$

Если бы в СММ было только одно состояние  $j$ , оценить параметр было бы легко. Оценки максимального правдоподобия величин  $\mu_j$  и  $\Sigma_j$  можно было бы получить простым усреднением:

$$\mu_j^* = \frac{1}{T} \sum_{t=1}^T o_t, \quad (8)$$

$$\Sigma_j^* = \frac{1}{T} \sum_{t=1}^T (o_t - \mu_j)(o_t - \mu_j)^T. \quad (9)$$

На практике, разумеется, имеется несколько состояний, и невозможно непосредственно привязать векторы наблюдений к отдельным состояниям, поскольку лежащие в основе последовательности состояний неизвестны. На начальном этапе необходимо осуществить некоторую приближенную привязку векторов к состояниям, тогда можно использовать уравнения (8) и (9) для получения нужных начальных значений параметров. Затем, с использованием описанного ниже алгоритма Витерби, находится наиболее правдоподобная последовательность состояний, векторы наблюдений заново привязываются к состояниям, после чего вновь используются уравнения (8) и (9) для получения лучших начальных значений. Этот процесс повторяется до тех пор, пока оценки перестают изменяться.

Так как полное правдоподобие каждой последовательности наблюдений основывается на суммировании всех возможных последовательностей состояний, каждый вектор наблюдения  $o_t$  вносит свой вклад в расчеты

значений параметров максимального правдоподобия для каждого состояния  $j$ . Иными словами, вместо того чтобы привязывать каждый вектор наблюдения к определенному состоянию, как это делалось в вышеупомянутом приближении, каждое наблюдение привязывается к каждому состоянию пропорционально вероятности состояния модели при наблюдении этого вектора. Таким образом, если  $L_j(t)$  обозначает вероятность пребывания в состоянии  $j$  в момент времени  $t$ , приведенные выше уравнения (8) и (9) становятся следующими взвешенными средними:

$$\mu_j^* = \frac{\sum_{t=1}^T L_j(t) o_t}{\sum_{t=1}^T L_j(t)}, \quad (10)$$

$$\Sigma_j^* = \frac{\sum_{t=1}^T L_j(t) (o_t - \mu_j)(o_t - \mu_j)'}{\sum_{t=1}^T L_j(t)}, \quad (11)$$

где суммирование в знаменателях обеспечивает требуемую нормализацию.

Уравнения (10) и (11) описывают процедуру рекуррентного оценивания Баума-Уэлша для средних и дисперсий НММ. Аналогичная, но несколько более сложная, процедура может быть получена для вероятностей перехода.

Конечно, чтобы применить соотношения (8) и (9), нужно рассчитать вероятность состояния  $L_j(t)$ . Это эффективно делается с использованием так называемого *алгоритма прямого-обратного хода* (*Forward-Backward algorithm*). Пусть прямая вероятность  $\alpha_j(t)$  для некоторой модели  $M$  с  $N$  состояниями определена в виде

$$\alpha_j(t) = P(o_1, \dots, o_t, x(t) = j / M). \quad (12)$$

Хотя это распределения плотностей, а не обычные вероятности, такое предположение удобно. Т.е.  $\alpha_j(t)$  — совместная вероятность наблюдения первых  $t$  векторов речи для состояния  $j$  в момент времени  $t$ . Эта прямая вероятность может быть эффективно рассчитана по следующей рекуррентной формуле:

$$\alpha_j(t) = \left[ \sum_{i=2}^{N-1} \alpha_i(t-1) \alpha_{ij} \right] b_j(o_t). \quad (13)$$

Вид этой рекуррентной формулы определяется тем обстоятельством, что вероятность пребывания в состоянии  $j$  в момент  $t$ , при наблюдениях  $o_t$ , можно вывести путем суммирования прямых вероятностей для всех возможных предшествующих состояний элемента состояний  $i$ , взвешенных вероятностями переходов  $a_{ij}$ . Несколько необычные пределы обусловлены тем обстоятельством, что состояния 1 и  $N$  не являются порождающими. Для понимания уравнений, содержащих непорождающие состояния в момент  $t$ , нужно полагать, что момент  $t - \delta t$  соответствует входному состоянию, а момент  $t + \delta t$  — выходному состоянию. Это важно тогда, когда модели связаны в последовательности так, что переходы через непорождающие

состояния происходят при переходах между моделями. Начальные условия для вышеупомянутого рекуррентного соотношения имеют вид

$$\alpha_1(1) = 1, \quad (14)$$

$$\alpha_j(1) = \alpha_1 b_j(o_1) \quad (15)$$

для  $1 < j < N$ , а конечные условия задаются так:

$$\alpha_N(T) = \sum_{i=2}^{N-1} \alpha_i(T) a_{iN}. \quad (16)$$

Заметим, что из определения  $\alpha_j(t)$  следует:

$$P(O/M) = \alpha_N(T). \quad (17)$$

Следовательно, вычисление прямой вероятности позволяет получить полное правдоподобие  $P(O/M)$ . Обратная вероятность  $\beta_j(t)$  определяется следующим образом:

$$\beta_j(t) = P(o_{t+1}, \dots, o_T / x(t) = j, M). \quad (18)$$

Как и в случае прямой вероятности, эта обратная вероятность может быть эффективно вычислена с использованием следующего рекуррентного соотношения:

$$\beta_i(t) = \sum_{j=2}^{N-1} a_{ij} b_j(o_{t+1}) \beta_j(t+1) \quad (19)$$

с начальным условием

$$\beta_i(T) = \alpha_{iN} \quad (20)$$

для  $1 < i < N$ , и конечным условием

$$\beta_1(1) = \sum_{j=2}^{N-1} a_{1j} b_j(o_1) \beta_j(1). \quad (21)$$

Заметим, что в приведенных выше определениях прямая вероятность является совместной вероятностью, тогда как обратная вероятность — условной. Это несколько асимметричное определение является преднамеренным, поскольку позволяет определить вероятность нахождения в состоянии как произведение этих двух вероятностей. По определению,

$$\alpha_j(t) \beta_j(t) = P(O, x(t) = j / M). \quad (22)$$

Отсюда

$$L_j(t) = P(x(t) = j / O, M) = \frac{P(O, x(t) = j / M)}{P(O/M)} = \frac{1}{P(O/M)} \alpha_j(t) \beta_j(t). \quad (23)$$

Из всего вышесказанного следует, что параметры СММ рекуррентно оцениваются по единственной последовательности наблюдения, т.е. по единственному экземпляру произнесенного слова. На практике же для

получения хороших оценок параметра, необходимо много экземпляров одного и того же слова. Тем не менее использование многократных последовательностей наблюдения не приводит к усложнению алгоритма. В заключение нужно упомянуть об одном моменте: вычисление прямых и обратных вероятностей связано с вычислением произведения большого количества вероятностей. На практике это означает, что числа становятся очень маленькими. Следовательно, во избежание проблем с малыми числами, вычисления реализуются с использованием логарифмирования.

**Алгоритм распознавания речи.** В предыдущем разделе описаны основные идеи рекуррентной оценки параметров СММ с использованием алгоритма Баума-Уэлша. При этом было отмечено, что эффективный рекурсивный алгоритм вычисления прямой вероятности позволяет заодно вычислить и полную вероятность  $P(O/M)$ . Таким образом, этот алгоритм может использоваться для нахождения модели, максимизирующей значение  $P(O/M_i)$ , и, следовательно, может использоваться для распознавания.

На практике, однако, удобнее осуществлять распознавание, основываясь на максимизации правдоподобия последовательности состояний, поскольку это легко обобщить на случай слитной речи, что невозможно при использовании полной вероятности. Это правдоподобие вычисляют, используя, в сущности, тот же алгоритм, что и при вычислении прямой вероятности, с тем лишь отличием, что суммирование заменяется поиском максимума. Предположим, что  $\phi_j(t)$  для данной модели  $M$  представляет максимальное правдоподобие наблюдения последовательности векторов речи от  $o_1$  до  $o_t$  и пребывания в состоянии  $j$  в момент времени  $t$ . Это частное правдоподобие можно эффективно вычислить с использованием следующего рекуррентного соотношения (сравните с (13)):

$$\phi_j(t) = \max_i \{\phi_i(t-1) a_{ij}\} b_j(o_t), \quad (24)$$

где

$$\phi_1(1) = 1, \quad (25)$$

$$\phi_j(1) = a_{1j} b_j(o_1) \quad (26)$$

для  $1 < j < N$ . Тогда максимальное правдоподобие  $P(O|M)$  имеет вид

$$\phi_N(T) = \max_i \{\phi_i(T) a_{iN}\}. \quad (27)$$

Что касается рекуррентного оценивания, прямое вычисление правдоподобия ведет к потере значащих разрядов, поэтому вместо него вычисляют логарифм правдоподобия. Таким образом, вместо уравнения (23) получаем

$$\psi_j(t) = \max_i \{\psi_i(t-1) + \log(a_{ij})\} + \log(b_j(o_t)). \quad (28)$$

Это рекуррентное соотношение является основой так называемого алгоритма Витерби.

**Неточности текстового сопровождения звукового материала.** При построении акустических моделей фонем предполагается, что последовательности фонем, построенные по текстовому описанию, в точности соответствуют звуковому сигналу.

Если в текстовом описании есть несоответствия, то процедура оценивания параметров моделей фонем сильно ошибается, поскольку оценка параметров происходит на звуковом материале, отнесенном к совсем другой фонеме. При этом ухудшается качество акустических моделей фонем, что приводит в дальнейшем к значительным ошибкам при распознавании речи.

Неточности в текстовом описании нового звукового материала можно разделить на три группы:

- 1) пропуск текста, описывающего звуковой сигнал;
- 2) добавление текста, который не произносился;
- 3) замена текста другим.

Третий случай можно рассматривать как комбинацию двух первых.

**Речевой и текстовый материал.** Для обучения акустической модели (АМ) в предварительном варианте системы распознавания был использован Акустический корпус украинской эфирной речи (АКУЕМ) [5]. В АКУЕМ представлена подготовленная и спонтанная речь, в основном прозвучавшая в украинском теле- и радиоэфире. Подготовленная речь в основном содержится в новостных радио- и телепередачах (чтение текста одним диктором), спонтанная — в различных ток-шоу (много дикторов, иногда перебивающих друг друга, эмоциональная спонтанная речь, речь с акцентом). Небольшая часть корпуса представляет собой аудиозаписи реальных судебных заседаний.

Общая длительность речи, сопровождаемой орфографической транскрипцией, — 500 часов, общее количество словоформ — около 1 100 000. Количество уникальных словоформ — 55 000. Часть корпуса АКУЕМ (82 часа) подробно аннотирована экспертами, в частности, отмечены акустические сегменты плохого качества (фоновый шум, музыка и речевые сбои (хезитации, оговорки и т.п.).

Речевой материал, использованный для построения предварительной АМ, состоял из аудиозаписей длительностью около 50 часов и содержал речь больше 2000 дикторов. Большинство дикторов представлено короткими записями, однако у 300 дикторов длительность записей составляет более 10 минут.

Для исследования возможности увеличения объема обучающего материала за счет неточно аннотированных речевых данных был проведен эксперимент на новом материале телеканала NewsOne. Этот материал состоял из 103 часов аудиозаписей и их текстовых подстрочников (сопроводительных текстов, не в точности соответствующих звуку). Для оценки количества неточностей в текстах небольшая часть текстового сопровождения была тщательно исправлена экспертами. Сравнение исправленного текста с исходным показало, что в исходном тексте содержится около 15 % ошибок.



Текстовый материал, использованный для обучения лингвистической модели, состоял из текстов корпуса АКУЕМ, текстов, взятых из Интернета, и нескольких искусственно созданных текстов.

**Система распознавания украинской речи.** В данной работе используется инструментарий НТК [6] на основе СММ. Этот инструментарий НТК использовался для построения акустических и лингвистических моделей, а для распознавания речи был разработан программный комплекс, который совместим с акустическими и лингвистическими моделями, построенными инструментарием НТК.

**Предварительная обработка речевого сигнала.** Речевой сигнал преобразуется в последовательность векторов признаков с интервалом анализа 25 мс и шагом анализа 10 мс. Вначале речевой сигнал фильтруется фильтром высоких частот с характеристикой  $P(z)=1-0,97z^{-1}$ . Затем применяется окно Хэмминга и вычисляется быстрое преобразование Фурье. Спектральные коэффициенты усредняются с использованием 26 треугольных окон, расположенных в мел-шкале, и вычисляются 12 кепстральных коэффициентов.

Логарифм энергии добавляется в качестве 13-го коэффициента. Эти 13 коэффициентов расширяются до 39-мерного вектора параметров путем дописывания первой и второй разностей векторов коэффициентов, соседних по времени. Для учета влияния акустического канала применяется вычитание среднего кепстра.

**Акустическая модель.** В качестве акустических моделей используются СММ. 56 украинских контекстно-независимых фонем (включая фонему-паузу) моделируются тремя состояниями марковской цепи без пропусков. Используется диагональный вид гауссовских функций плотности вероятности.

Редко встречающиеся фонемы моделируются 64 смесями гауссовских функций плотности вероятности, более часто встречающиеся фонемы моделируются большим числом смесей, наиболее часто встречающиеся фонемы используют 1024 смесей.

Словарь транскрипций создается автоматически из орфографического словаря в 125 000 словоформ с использованием правил преобразования буква-фонема.

**Обучение акустической модели при неполной информации.** Использовался аудио- и текстовый материал телеканала NewsOne [7]. Этот материал был разделен на обучающую (100 часов) и тестовую (3 часа) части. Исходные файлы представляют собой новостные сообщения длительностью от 1 до 10 минут. В общем случае удобнее работать с относительно короткими речевыми сегментами, чтобы локализовать возможные ошибки определения границ фонем.

Обучение АМ на новом речевом материале состояло из следующей последовательности операций над каждым файлом:

1. Составление списка слов, входящих в файл, и автоматическое порождение фонемных транскрипций для этих слов.
2. Обучение биграммной ЛМ.

3. Разбиение речевых файлов на сегменты.
4. Распознавание речевых сегментов предварительной системой распознавания для получения автоматических транскрипций.
5. Выравнивание подстрочника с автоматически полученными транскрипциями с использованием алгоритма динамического программирования для выявления речевых сегментов, где есть точное соответствие звука и текста.

В результате выполнения данной процедуры для всех файлов обучающей выборки было выявлено 60 часов речи, где есть точное соответствие текстового сопровождения и звукового сигнала. Данный материал использовался для обучения АМ и построения новой системы распознавания.

**Экспериментальные результаты.** Эксперименты по распознаванию речи проводились для тестовой выборки длительностью в 3 часа. Для сравнения исследовалась точность распознавания речи для предварительной системы распознавания, построенной на акустическом материале АКУЕМ, размеченном экспертами. Пословная ошибка распознавания составила 26,56 %. Для новой системы распознавания, построенной на акустическом материале, размеченном автоматически, пословная ошибка составила 19,97 %. Таким образом, разработанная технология позволила снизить процент ошибок на 24,8 %.

Следует заметить, что в тестовой выборке около 10 % файлов было зашумлено различными неречевыми сигналами (в основном музыкой). Для этих файлов был получен значительно больший процент ошибок распознавания — больше 30 %. Эти сегменты легко выявить при помощи автоматической процедуры оценивания отношения сигнал/шум.

**Выводы.** Разработанная технология построения акустических моделей использует неполное текстовое описание звукового материала. Особенностью технологии является автоматическое выявление фрагментов звукового сигнала, где есть точное соответствие текстового описания и звукового сигнала. Данный подход может применяться для автоматической разметки большого количества частично аннотированных звуковых сигналов, что значительно уменьшает затраты на разработку систем распознавания речи. Экспериментальные результаты показывают эффективность использованного подхода и уменьшение количества ошибок при распознавании речи на 24,8 %.

1. *Lamel L., Gauvain J.L., Adda G.* Lightly Supervised Acoustic Model Training // ISCA ITRW Workshop on Automatic Speech Recognition: Challenges for the new Millenium. — Paris, 2000. — P. 150–154.
2. *Lamel L., Gauvain J., Adda G.* Investigating lightly supervised acoustic model training // Proceedings IEEEICASSP. — 2001. — Vol. 1. — P. 477–480.
3. *Lamel L., Gauvain J.L., Adda G.* Lightly supervised and unsupervised acoustic model training // Computer Speech and Language. — 2002. — 16(1). — P. 115–229.
4. *Paulik M., Waibel A.* Lightly supervised acoustic model training on EPPS recordings // INTERSPEECH. — 2008. — P. 224–227.
5. Ukrainian Broadcast Speech Corpus Development / V. Pylypenko, V. Robeiko, M. Sazhok et al. // Proceedings of SPECOM 2011, 14-th International Conference “Speech and Computer”. Kazan, Russia, 2011. — Kazan: Paladin, 2011. — P. 435–440.

6. The HTK Book Version 3.4.1 / S.J. Young, G. Evermann, M.J.F. Gales et al. — Cambridge: Cambridge University, 2009. — 384 p.
7. <http://newsone.ua/>

Международный научно-учебный центр  
информационных технологий и систем  
НАН Украины и Министерства образования  
и науки, молодежи и спорта Украины, Киев

Получено 04.10.2012