

# Гибридная технология перевода

Юлия Епифанцева  
PROMT



**PROMT**<sup>®</sup>



# Машинный перевод

Машинный (автоматический) перевод – процесс перевода текстов с одного естественного языка на другой с помощью компьютерной программы

# Основные типы систем МП

- Rule-based машинный перевод (RBMT) – перевод, основанный на правилах.
  - Статистический машинный перевод (SMT).
- ➔ Гибридные системы перевода (HMT)

# Типы RBMT

- **Системы по типу Transfer**

предложение на языке входа =>

морфологический, грамматический, семантический анализ =>

преобразование в структуру выходного языка (TRANSFER) =>

синтез выходного предложения по полученной структуре=>

предложение на языке выхода

- **Системы по типу Interlingua**

предложение на языке входа =>

анализ входного предложения в терминах метаязыка =>

синтез из метаструктуры предложения выходного языка =>

предложение на языке выхода

Разработка метаязыка = языконезависимое представление, наличие знаний о мире (онтологии, логики предикатов)

# Компоненты RBMT на примере PROMT

- Лингвистические базы данных
  - двуязычные словари
  - файлы имен, транслитерации
  - морфологические таблицы
- Модуль перевода
  - грамматические правила
  - алгоритмы перевода

# Двуязычные словари

имеют трехуровневую структуру для настройки системы на различные предметные области:

- **Генеральный словарь** (от 50 до 250 тысяч статей)
- **Специализированные словари** (от 5 до 100 тысяч статей; охватывают различные тематики: бизнес, спорт, IT, добыча нефти и газа, металлургия...)
- **Пользовательские словари** (вспомогательные, открыты для редактирования пользователю)

# Словарная статья (1)

- Слово или выражение на входном языке
- Коллекция переводов (активных и неактивных)
- Структурированное описание различной лингвистической информации (морфологической, синтаксической, семантической) в виде набора признаков и модификаторов

# Словарная статья (2)



perfectly shaped a  
perfectly well d  
perfecto n  
perfervid a  
perficient a  
perfidе a  
perfidious a  
perfidy n  
perforate v  
perforating n  
perforation n  
perforative a  
perforator n  
perforce d  
perform v  
perform on stage v  
perform [...] about-face v  
perform [...] act v  
perform [...] surgery v  
**performance n**  
performance art n  
performance bond n  
performance car n  
performance characteristic n  
performance chart n

**performance** *Существительное*

- 1. работа *f*
- 2. выполнение *n*
- 3. действия *pl* (activity)
- 4. выступление *n* (theatre)
- 5. показатели *pl*
- 6. результаты *pl*
- 7. представление *n* (theater)
- 8. исполнение *n* (performance of smth)
- 9. выступление *n* (action)
- 10. выступление *n* (during performance)
- 11. показатели *pl*
- 12. работа *f*
- 13. поведение *n*
- 14. исполнение *n* (his performance of smth)
- 15. оценки *pl*
- 16. успеваемость *f* (of students)
- 17. результаты *pl*
- 18. работа *f* (brake performance)
- 19. показатели *pl*
- 20. успехи *pl* (performance of a country or company)
- 21. выполнение *n*
- 22. работа *f* (brake performance)
- 23. производительность *f* (vector performance of smth)
- 24. исполнение *n* (performance of dance)



# Словарная статья (3)

Словарная настройка на примере 'chest'

| Модификаторы      | Перевод              |
|-------------------|----------------------|
| гXj               | груд/н13/1           |
| гX_H8             | комод/н1/1           |
| гPpb*Yj           | груд/н13/1           |
| гY_H8*Amn         | комод/н1/1           |
| гPцю_ia_IF_IE_IG* | фонд/н1/1            |
| гVhOyv*Yr         | упаков/н4/1          |
| FY                | грудной клетки/а98/1 |
| Fx                | грудн/а5/1           |
| гPtv*Y_H8         | комод/н1/1           |
| гP_Rn*Y_H8        | комод/н1/1           |

| Контекст  | Перевод   |
|---|---|
| This is a chest   | Это грудь (комод)   |
| There is a tattoo on his chest  | На его груди есть татуировка  |
| At the time of the chest examination the blood pressure may be taken  | Во время обследования грудной клетки может быть измерено давление   |
| The University Chest is a term used at Oxford in connection with the financial aspects of the university and its administration | Университетский фонд – термин, использованный в Оксфорде в связи с финансовыми аспектами университета и его администрации |
| The oak chest with iron lock  | Комод из дуба с железным замком   |

# Уровни анализа предложения (1)

- **Препроцессор**
- **Нормализация текста** (удаление повторяющихся пробелов...)
- **Токенизация входной цепочки** (поиск входных словоформ в словаре с сопутствующим морфологическим анализом)
- **Лексический анализ** (контекстный анализ, различные склейки: имена, номера телефонов, даты...)
- **Снятие омонимии** (определение частей речи в случаях грамматической неоднозначности)
- **Уровень сбора групп** (соединение лексических единиц в группы)
- **Анализ сложных предложений** (выделение простых в составе сложного, синтаксическая омонимия)
- **Семантико-синтаксический разбор** (заполнение глагольного фрейма)
- **СИНТЕЗ** (синтез по полученной структуре, расстановка элементов внутри группы и групп в предложении...)

# Уровни анализа предложения (2)

French restaurants and bars, Mr. Felise notes, are getting more popular in the USA.

- Лексический анализ

French restaurants and bars , [Mr. Felise ] notes , are getting [ more popular ] in the USA .

- Синтаксические группы

[ French restaurants ] and [ bars ] , [ Mr. Felise ] notes [ , ] , [ are getting ] more popular [ in ] the USA [ . ]

- Анализ сложного предложения

{ [ **French restaurants** ] and [ bars ] { [ , ] [ **Mr. Felise** ] notes [ , ] , } [ are getting ] more popular [ in ] the USA [ . ] }

- Анализ простых предложений

{ [ **French restaurants** ] and [ bars ] { [ , ] [ **Mr. Felise** ] notes [ , ] , } [ are getting ] more popular [ in ] the USA [ . ] }

# Уровни анализа предложения (3)

```
{S
  {S
    {NP-SBJ(JJ French)(NNS restaurants)}
    (CC and)
    {NP-SBJ(NNS bars)}
    {S
      (.)
      {NP-SBJ(NNP (Mr. Felise))}
      {VP
        {VP(VBx notes)}
        (.)}}
    {VP
      {VP(- are)
        {VP(VBx getting)}(- more)
        {AdjP(JJR (more popular))
          {PP(IN in)
            {NP(DT the)(NNP USA)}}}}
      (.)}}
```

# Преимущества и недостатки RBMT

- **Преимущества**
  - синтаксическая и морфологическая точность,
  - стабильность и предсказуемость результата,
  - возможность настройки на предметную область.
- **Недостатки**
  - трудоемкость и длительность разработки,
  - необходимость поддерживать и актуализировать лингвистические БД,
  - «машинный акцент» при переводе.

# Статистический машинный перевод (1)

## История

- Принципы SMT разработаны еще в 1949 г. Уорреном Уивером
- «Вторая волна» – начало 1990-х, IBM
- «Третья волна» – Google, Microsoft, Language Weaver, Яндекс и десятки других

Статистический МП – это поиск наиболее вероятного перевода предложения с использованием данных, полученных из параллельных корпусов.

# Статистический машинный перевод (2)

- Сегодня SMT – наиболее активно разрабатываемая архитектура МТ. Почему?
  - Легко построить, если есть двуязычный корпус, ноль / минимум лингвистики
  - Переносимость технологии на любые пары языков
  - Лексическая гладкость
- Ограничения / недостатки:
  - Ограниченность параллельных корпусов в природе и их real-life качество
  - Плохо справляется с морфологией / синтаксисом (по сравнению с RBMT)
  - Искажение информации (дублирование, пропуск или подмена информации)

**USA is to blame = США не виноват**

**Russia is to blame = Россия виновата**

# Выводы

Обе технологии имеют свои достоинства и недостатки, но главное – они не решили задачу по получению качественного машинного перевода.

MT-сообщество ожидает прорыва в качестве перевода в гибридных моделях RBMT + SMT.



# Гибридная технология PROMT

Объединение RBMT и статистических технологий

- лингвистический анализ входного предложения
- порождение вариантов перевода
- использование статистических технологий
- оценка и выбор лучшего варианта перевода с использованием Модели языка

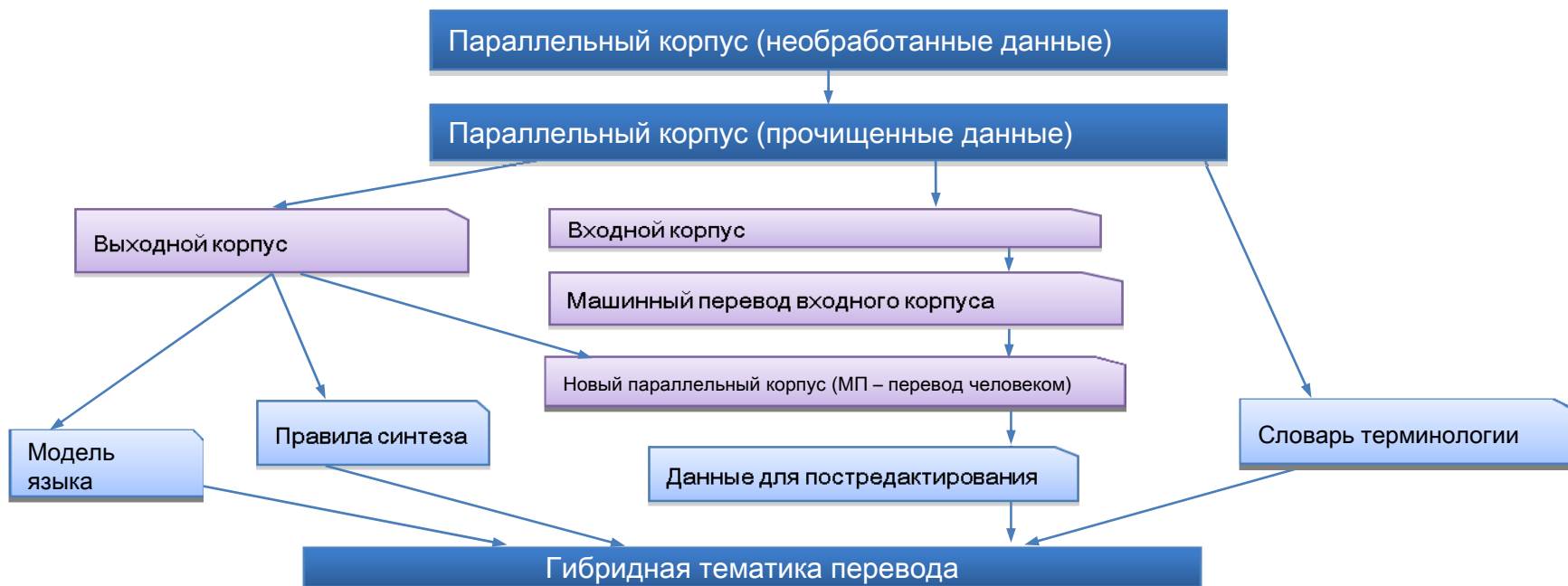
# Этапы Гибридной технологии

- Обучение RBMT на основе параллельного корпуса с использованием статистических технологий
- Эксплуатация на основе натренированной системы

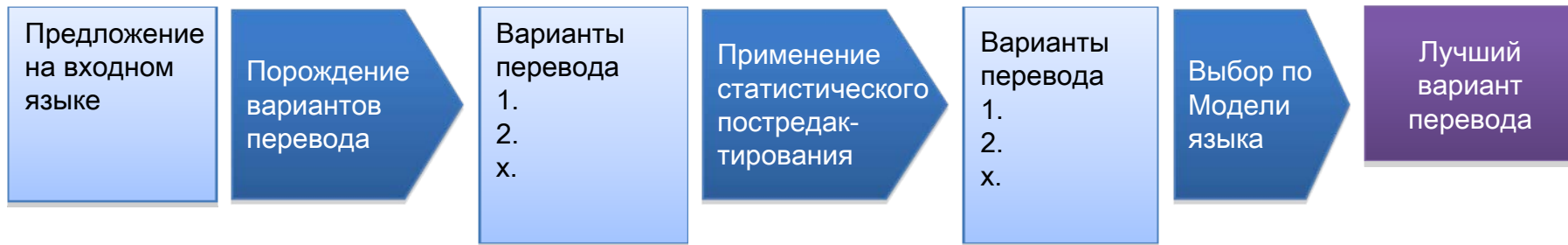
# Архитектура Гибридной технологии



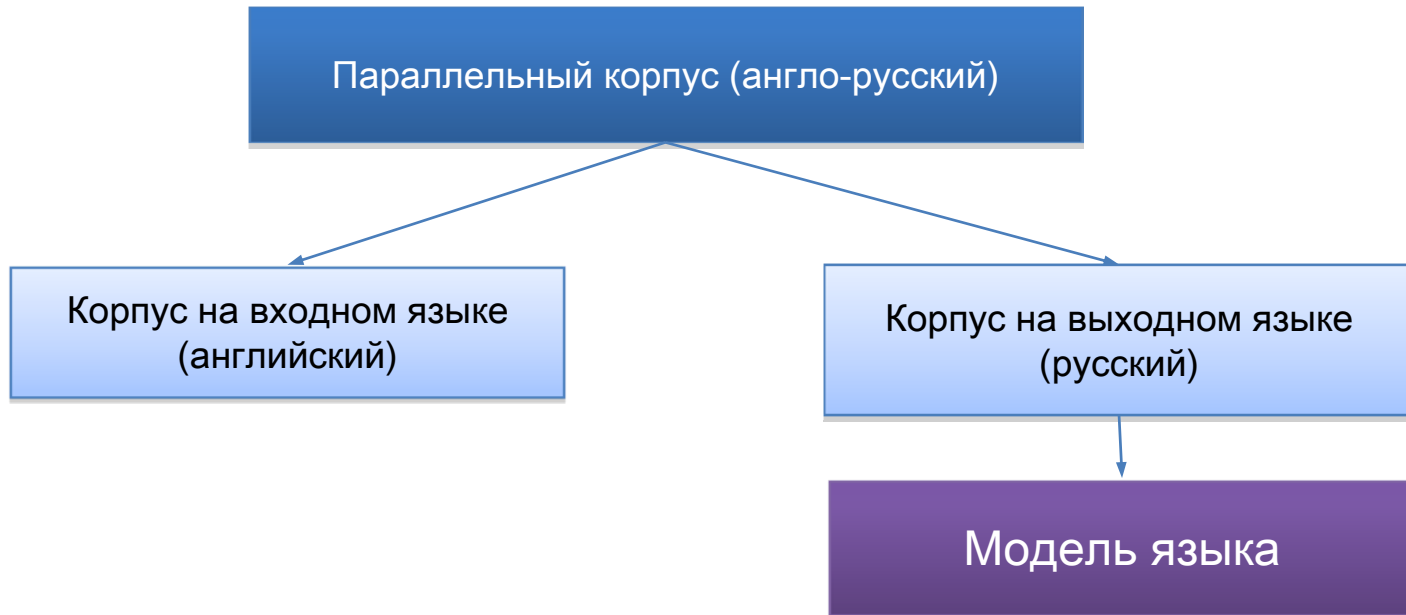
# Обучение



# Эксплуатация



# Модель языка (1)



# Модель языка (2)

- Модель языка (Language Model/LM) – это набор  $n$ -грамм монопольного корпуса с их вероятностными характеристиками.
- $N$ -грамма – это последовательность слов из предложений длины  $n$ .

# Модель языка (3)

|            |                                     |             |
|------------|-------------------------------------|-------------|
| -1.755048  | по двум очередям в                  | -0.6624343  |
| -1.755048  | не к очередям каждого               | -0.0527804  |
| -1.755048  | относящимся к очередям относятся    | -0.05122143 |
| -1.755048  | обслуживания различным очередям и   | -0.05385702 |
| -1.056816  | обеими входящими очередями и        | -0.6624451  |
| -0.4883884 | между двумя очередями </s>          |             |
| -1.056816  | совместно используется очередями и  | -0.6624537  |
| -2.103502  | буфера между очередями путем        | -0.05193719 |
| -2.103502  | между различными очередями процент  | -0.04642604 |
| -2.188747  | вызовы в очередях по                | -0.05341507 |
| -1.317873  | использовать в очередях виртуальных | -0.5909631  |
| -1.044643  | находящимся в очередях уровня       | -0.05364625 |
| -1.044643  | пакеты в очередях уровня            | -0.05345588 |
| -2.188747  | потерями в очередях входа           | -0.6321825  |
| -2.118226  | Во входящих очередях SRR            | -0.6562098  |
| -2.118226  | и выходной очередях количество      | -0.0537526  |
| -0.435618  | в других очередях </s>              |             |
| -1.55669   | особенно при очередях с             | -0.05344867 |
| -1.641104  | снятия изоляции очистите от         | -0.662367   |
| -1.641104  | головного узла очистите статистику  | -0.05261082 |
| -0.8563652 | <s> Чтобы очистить базу             | 0.1404502   |
| -2.303862  | поздних версиях очистить ULAN       | -0.1768938  |
| -2.303862  | увеличится если очистить место      | -0.1765835  |
| -0.435618  | NAT можно очистить с                | -0.6248463  |
| -1.742327  | Если необходимо очистить сеанс      | -0.6620195  |
| -2.303862  | если необходимо очистить все        | -0.05370406 |



# Модель языка (4)

- Perplexity (PPL) – вычисляемая для предложения «степень его искаженности» по отношению к данной LM. Чем меньше PPL, тем «естественнее» предложение.
- Модель языка
  - оценка релевантности (через PPL) каждого перевода по отношению к данному корпусу,
  - выбор лучшего варианта среди всех порожденных.

# Как работает Гибридная технология

- Создание терминологического словаря из параллельных текстов для RBMT автоматическим путем
- Порождение всех возможных вариантов перевода на основе
  - лексических вариантов
  - вариантов синтеза разных конструкций
  - применения постредактирования
- ➔ выбор лучшего варианта через Модель языка

# Терминологический словарь (1)

Технология получения:

а) на основе параллельного корпуса составляются таблицы  $n$ -грамм входного корпуса вместе с вероятностями соответствий этих  $n$ -грамм  $n$ -граммам выходного корпуса,

б) на основании парсеров для входного и выходного языков, а также частотных характеристик из общего числа  $n$ -грамм извлекаются релевантные для словаря единицы с некоторой грамматической информацией (например, о валентности)

→ создается двуязычный глоссарий

в) в автоматическом режиме создается словарь для RBMT

# Терминологический словарь (2)

|    | А                            | В                              |
|----|------------------------------|--------------------------------|
| 1  | access control change        | изменение прав доступа         |
| 2  | access control permission    | разрешение прав доступа        |
| 3  | access control policy        | политика прав доступа          |
| 4  | access detail                | подробная информация о доступе |
| 5  | access information           | информация о доступе           |
| 6  | access key generation        | создание ключей доступа        |
| 7  | access policy                | политика доступа               |
| 8  | access policy rule           | правило политики доступа       |
| 9  | access rule                  | правило доступа                |
| 10 | acme organization            | организация ACM                |
| 11 | action button                | кнопка действий                |
| 12 | action icon                  | значок действий                |
| 13 | action item association      | связь распоряжен               |
| 14 | action item attribute        | атрибут распоряж               |
| 15 | action item data             | данные распоряж                |
| 16 | action item information page | страница свойств               |

activity chart n  
activity completion n  
activity deadline n  
activity due date n  
activity icon n

**activity deadline**, *Существительное*

- 1. наступление крайнего срока выполнения задачи *л*

# Лексические варианты

The restaurant features a unique **space** with a cozy lounge and an eclectic blend of music, art and sculpture.

## Rule-based

Ресторан представляет собой уникальное **пространство (место)** с удобным залом и эклектичной смесью музыки, искусства и скульптуры.

## Hybrid

Ресторан представляет собой уникальное **пространство** с удобным залом и эклектичной смесью музыки, искусства и скульптуры. (PPL==778)

Ресторан представляет собой уникальное **место** с удобным залом и эклектичной смесью музыки, искусства и скульптуры. (PPL=522)

# Варианты синтеза конструкций (1)

Rule-based: выбор определенной модели синтеза

Hybrid: синтезирование нескольких вариантов перевода

Правило синтеза: перевод конструкции to + инфинитив

- чтобы + инфинитив
- для + существительное

You can use the same steps **to edit the collection**.

Можно использовать те же самые шаги, **чтобы отредактировать коллекцию**. (*PPL=372*)

Можно использовать те же самые шаги **для редактирования коллекции**. (*PPL=358*)

# Варианты синтеза конструкций (2)

Rule-based: выбор определенной модели синтеза

Hybrid: синтезирование нескольких вариантов перевода

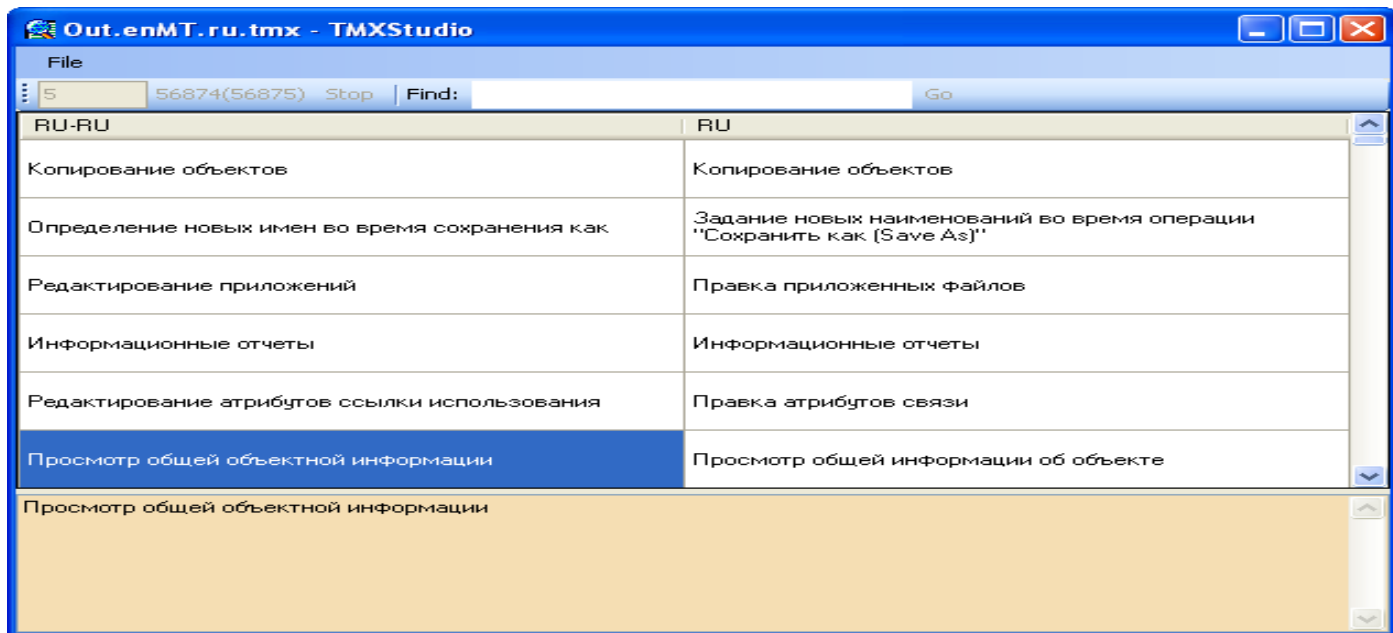
Правило синтеза: порядок следования подлежащего и сказуемого.

Click Browse to browse the path for the folder in which you want **newly created documents to be located**.

Нажмите Browse для просмотра пути к папке, в которой Вы хотите, чтобы **недавно созданные документы были расположены**. (PPL= 290)

Нажмите Browse для просмотра пути к папке, в которой Вы хотите, чтобы **были расположены недавно созданные документы**. (PPL= 274)

# Данные для постредактирования (1)



The screenshot shows the TMXStudio interface with a table comparing source (RU-RU) and target (RU) text for post-editing. The table has two columns: 'RU-RU' and 'RU'. The first row is highlighted in blue. Below the table, there is a section for 'Просмотр общей объектной информации'.

| RU-RU  | RU   |
|--|--|
| Копирование объектов                           | Копирование объектов   |
| Определение новых имен во время сохранения как | Задание новых наименований во время операции "Сохранить как (Save As)" |
| Редактирование приложений                      | Правка приложенных файлов  |
| Информационные отчеты                          | Информационные отчеты  |
| Редактирование атрибутов ссылки использования  | Правка атрибутов связи   |
| Просмотр общей объектной информации            | Просмотр общей информации об объекте                                   |

Просмотр общей объектной информации



# Данные для постредактирования (2)

Технология : на основе параллельного корпуса выделяется таблица n-грамм входного корпуса вместе с вероятностями соответствий этих n-грамм n-граммам выходного корпуса.

- *с платежом PayPal банковским переводом → в случае платежа PayPal посредством банковского перевода*
- *вводите банковский перевод → инициируете перевод*
- *когда Вы закончены → после окончания Вашей работы*
- *каждое усилие было приложено → были предприняты все усилия*

# Данные для постредактирования (3)

Пример применения нескольких замен сегментов машинного переводами сегментами человеческого перевода.

With PayPal payment by bank transfer, you initiate a bank transfer from your bank account to your PayPal account.

С платежом PayPal банковским переводом вы вводите банковский перевод с Вашего банковского счета на ваш счет PayPal. (PPL=95)

В случае платежа PayPal посредством банковского перевода вы инициируете перевод с Вашего банковского счета на ваш счет PayPal. (PPL == 7)

Исходный текст

*Click Browse to browse the path for the folder in which you want newly created documents to be located.*

Порождение лексических вариантов

*Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы недавно созданные документы были **расположены**.  
Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы недавно созданные документы были **размещены**.*

Порождение вариантов синтеза

*Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы недавно созданные документы **были расположены**.  
Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы недавно созданные документы **были размещены**.  
Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы **были расположены** недавно созданные документы.  
Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы **были размещены** недавно созданные документы .*

Порождение вариантов постредактирования

*Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы недавно созданные документы **были расположены**.  
Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы недавно созданные документы **были размещены**.  
Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы **были расположены** недавно созданные документы .  
Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы **были размещены** недавно созданные документы .  
Нажмите Browse **для просмотра** пути к папке, в которой Вы хотите, чтобы недавно созданные документы **были расположены**.  
Нажмите Browse **для просмотра** пути к папке, в которой Вы хотите, чтобы недавно созданные документы **были размещены**.  
Нажмите Browse **для просмотра** пути к папке, в которой Вы хотите, чтобы **были расположены** недавно созданные документы.  
Нажмите Browse **для просмотра** пути к папке, в которой Вы хотите, чтобы **были размещены** недавно созданные документы .*

```
graph TD; A[ ] --> B[Оценка LM]; B --> C[Выбор лучшего варианта];
```

Оценка LM

Выбор лучшего варианта

*Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы недавно созданные документы были расположены. (PPI = 556)*

*Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы недавно созданные документы были размещены. (PPI = 601)*

*Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы были расположены недавно созданные документы. (PPI = 526)*

*Нажмите Browse, чтобы рассмотреть путь к папке, в которой Вы хотите, чтобы были размещены недавно созданные документы. (PPI = 569)*

*Нажмите Browse для просмотра пути к папке, в которой Вы хотите, чтобы недавно созданные документы были расположены. (PPI = 277)*

*Нажмите Browse для просмотра пути к папке, в которой Вы хотите, чтобы недавно созданные документы были размещены. (PPI = 301)*

*Нажмите Browse для просмотра пути к папке, в которой Вы хотите, чтобы были расположены недавно созданные документы. (PPI = 261)*

*Нажмите Browse для просмотра пути к папке, в которой Вы хотите, чтобы были размещены недавно созданные документы. (PPI = 283)*

*Нажмите Browse для просмотра пути к папке, в которой Вы хотите, чтобы были расположены недавно созданные документы. (PPI = 261)*

# LM Statistics



# Поиск по параллельным корпусам

Найдено: **178** вхождений, **3** документов  
Размер подкорпуса: **100%** от общего объема корпуса EMC(en)

**Английский: EMC(en)**      Русский : EMC(ru)

|   |   |
|---|---|
| On the server, <b>open</b> the Navisphere Server Utility.   | На сервере сервере откройте утилиту Navisphere Server Utility.  |
| <b>Open</b> a web browser and type your device's unique web address, which can be found on the Remote Access Settings page of the | Откройте веб-браузер и введите уникальный веб-адрес устройства, который можно найти на странице Удаленный доступ      |
| <b>Open</b> Replicator for Symmetrix Implementation   | Внедрение решения Open Replicator для Symmetrix   |
| When your device finishes rebooting, <b>open</b> the Manager.   | По завершении перезагрузки устройства откройте Диспетчер.   |
| Choose whether to display the standard action buttons for the tab (Edit, <b>Open</b> , View, and Export).                         | Выберите необходимость отображения стандартных кнопок действий для вкладки (Правка, <b>Открыть</b> , Вид, и Экспорт). |

Язык: english

Настроить корпус

словоформа      лемма

open      1

Омонимия:

До снятия

Часть речи

v

Синтаксис...

Пунктуация:

слева      справа

Позиция в предложении:

любая

Регистр:

любой

Искать отрицание

Настройки поиска: дополнительно ▾

Искать      Очистить

# Выводы

- **Преимущества RBMT сохраняются:**
  - синтаксическая и морфологическая точность,
  - стабильность и предсказуемость результата,
  - возможность настройки на предметную область.
- **Недостатки RBMT нивелируются за счет использования параллельных корпусов и статистических методов**
  - автоматическая настройка лингвистических баз данных (быстрое и качественное извлечение терминологии),
  - исчезает «машинный» акцент при перевода (варианты синтеза и постредактирование).



# Примеры перевода с помощью гибридной технологии

| Исходный текст   | Гибридный машинный перевод   | Перевод человека   |
|--|--|--|
| To replace any parts of devices or cables, use only original spare parts or parts which are explicitly licensed by the manufacturer. | Для замены любых <b>частей</b> устройств или кабелей <b>используйте только первоначальные запасные части или части</b> , явно лицензированные изготовителем.                             | Для замены любых <b>компонентов</b> устройств или кабелей <b>должны использоваться только оригинальные запасные части или компоненты</b> , явно лицензированные производителем.                      |
| Remove the faulty CF card from the SC card.  | <b>Удалите</b> неисправную CF-карту из <b>Платы</b> SC.  | <b>Выньте</b> неисправную CF-карту из <b>платы</b> SC.   |
| It is a precondition for reuse and recycling of used electrical and electronic equipment.  | Это является предварительным <b>условие</b> для вторичного использования и <b>сброса используемого</b> электрического и электронного оборудования..                                      | Это является <b>необходимым</b> предварительным условием для вторичного использования и <b>переработки отработавшего</b> электрического и электронного оборудования.                                 |
| Total Loss of Capacity Reception   | Total Loss of Capacity Transmit; Полная потеря приема пропускной способности   | Total Loss of Capacity Reception; Полная потеря пропускной способности при приеме  |
| Document history   | Хронология <b>документации</b>   | Хронология <b>документа</b>  |
| There is a risk of explosion if the battery is not replaced correctly!   | <b>Если батарея не заменяется правильно</b> , существует <b>риск взрыва!</b>   | При <b>неправильной</b> замене батареи существует <b>опасность взрыва!</b>   |
| This document is intended for anyone tasked with installing, operating, maintaining, or troubleshooting the SURPASS hiT 7300 system. | Этот документ предназначен для <b>специалистов</b> , для <b>которого определяют задачу</b> с инсталляцией, эксплуатацией, обслуживанием, или устранению отказов SURPASS hiT 7300 систем. | Этот документ предназначен для <b>персонала</b> , <b>отвечающего за процессы</b> установки, эксплуатации, технического обслуживания или устранения отказов применительно к системе SURPASS hiT 7300. |



# Примеры перевода с помощью гибридной технологии

| Исходный текст  | Гибридный машинный перевод  | Перевод человека  |
|---|---|---|
| Unfortunately, if access to your account remains limited for an extended period of time, it may result in further limitations or eventual account closure.            | К сожалению, <b>если доступ к счету остается ограниченным в течение длительного периода времени, он</b> может привести к дальнейшим ограничениям или <b>конечному закрытию счета.</b> | К сожалению, ограничение <b>доступа к учетной записи на более длительный период времени</b> может привести к дальнейшим ограничениям или к <b>ее окончательному закрытию.</b> |
| We thank you for your prompt attention to this matter.  | Мы благодарим Вас за <b>вашу оперативность по отношению к данному вопросу.</b>  | Благодарим вас за то, что вы сразу же обратили <b>внимание на эту проблему.</b>   |
| Your account was temporarily limited as a security measure and full access has not been restored yet because we haven't heard from you.                               | Ваш счет <b>был временно ограничен как меры безопасности, и полный доступ еще не был восстановлен, поскольку мы не получили известие от вас.</b>                                      | Ваш счет временно ограничен <b>с целью защиты. Полный доступ не был восстановлен, потому что мы не получили от вас ответа.</b>  |
| We understand it's frustrating and inconvenient to have your account permanently limited but you'll still be able to see your transaction history for a limited time. | Мы понимаем ваше разочарование и неудобства, связанные с постоянным ограничением <b>счета постоянно, но вы по-прежнему сможете просматривать историю транзакций.</b>                  | Мы понимаем ваше разочарование и неудобства, связанные с постоянным ограничением <b>счета, но некоторое время вы сможете видеть историю ваших транзакций.</b>                 |
| Email address provided is not a valid email!  | <b>Предоставленный адрес электронной почты не является допустимым электронным письмом!</b>  | <b>Указан недействительный адрес электронной почты!</b>   |

# Что дает гибридная технология перевода?

- быструю автоматическую настройку на основе translation memories заказчика,
- терминологическую точность перевода, а также единство стиля,
- получение дополнительных полезных данных – двуязычного терминологического словаря

*Разница в качестве между гибридным машинным переводом и переводом, выполненным человеком, составляет 25-40 % (== объем постредактирования)*

# Объемы данных

Минимальный объем данных для тренировки гибридной технологии:

- 50 000 сегментов из ТМ,
- не менее 500 000 слов в исходном языке.



Спасибо!

[www.promt.ru](http://www.promt.ru)