

Развитие машинного перевода

Балышев Артем

Машина Троянского

В 1933 году Советский ученый Пётр Троянский обращается в Академию Наук СССР с изобретённой им «машиной для подбора и печатания слов при переводе с одного языка на другой». Машина была крайне проста: большой стол, печатная машинка с лентой и плёночный фотоаппарат. На столе лежали карточки со словами и их переводами на четырёх языках. Машина Троянского впервые на практике реализовала тот самый «промежуточный язык» (interlingua).



Я I ICH YO	ХОТЕТЬ WANT WOLLEN QUERER	МНОГО MANY VIEL MUCHO	ХУРМА PERSIMMON PERSIMONE CACU
МЕСТ., ЕД. Ч., ИМ. П.	ГЛАГ., I. Л., ЕД. Ч., НЕСОВ., НАСТ. ВР., ДЕЙСТВ. ЗАЛОГ	ЧИСЛ., ИМ. П.	СУЩ., МН. Ч., РОД. П., НЕОДУШ.

КРАТКАЯ ИСТОРИЯ МАШИННОГО ПЕРЕВОДА

РВМТ

ПЕРЕВОД НА ОСНОВЕ ПРАВИЛ

ПО СЛОВАМ

ТРАНСФЕР-
НЫЙ

ИНТЕР-
ЛИНГВА

ЕВМТ

ПЕРЕВОД
НА
ПРИМЕРАХ

SMT

СТАТИСТИЧЕСКИЙ ПЕРЕВОД

НА СЛОВАХ

СИНТАК-
СИЧЕСКИЙ

НА ОСНОВЕ ФРАЗ

NMT

НЕЙРО-
СЕТЕВОЙ
ПЕРЕВОД

RNN
LSTM



Машинный перевод на основе правил — Rule-based Machine Translation (R^{DM}IT)



Среди плюсов RBMT отмечают морфологическую точность (не путает слова), воспроизводимость (все переводчики получают одинаковый результат) и возможность обучить специальным терминам под предметную область.

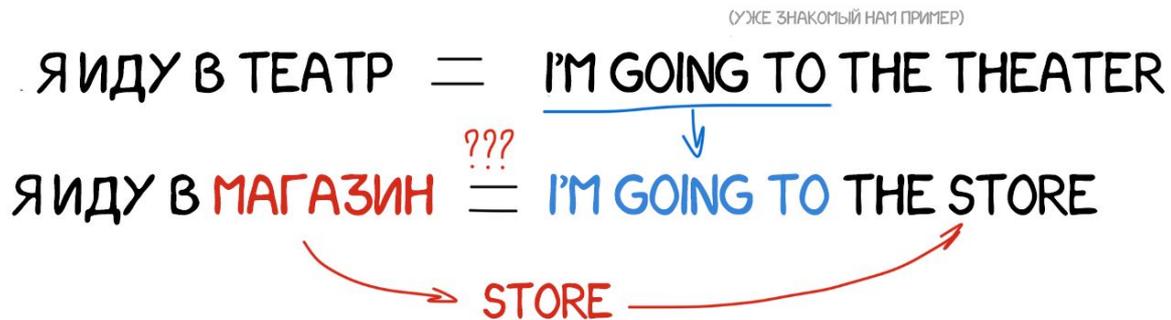
Минусы:

- исключения из правил языка - неправильные глаголы в английском, плавающие приставки в немецком, суффиксы в русском
- омонимия. Одно и то же слово может иметь разный смысл зависимости от контекста, а значит отличается и его перевод.

я	ХОЧУ	СОРОК	КИЛОГРАММ	ХУРМЫ
↓	↓	↓	↓	↓
I	WANT	FORTY	KILOGRAM	PERSIMMONS

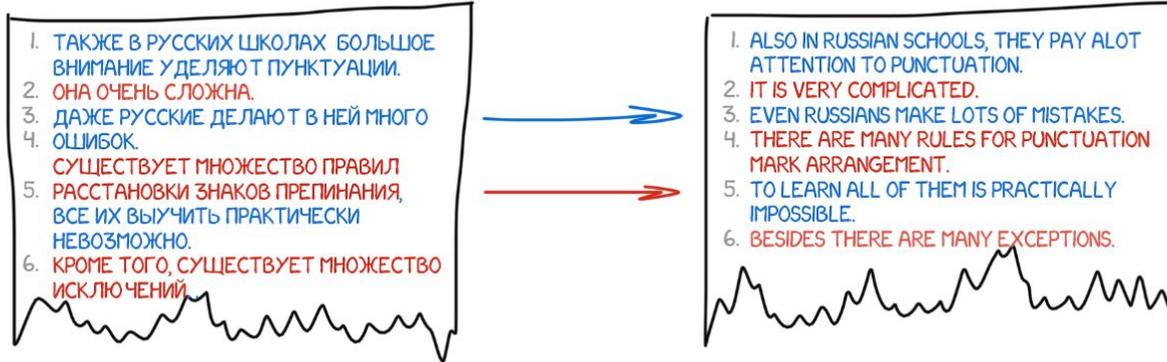
Машинный перевод на примерах — Example-based Machine Translation (EBMT)

Основной принцип: А что если не пытаться каждый раз переводить заново, а использовать уже готовые фразы?

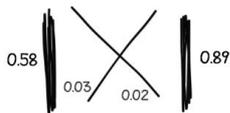


Статистический машинный перевод — Statistical Machine Translation (SMT)

ПАРАЛЛЕЛЬНЫЙ КОРПУС



THE HOUSE



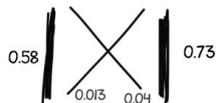
DAS HAUS

BLUE HOUSE



BLAUES HAUS

THE CAR



DAS AUTO

английский немецкий русский Перевести

leave ✓



Предложить исправление

уходить: варианты перевода

глагол

go	идти, ехать, ходить, ездить, уходить, проходить
go away	уходить, отойти, разойтись
depart	отступать, отклоняться, уходить, отправляться, уезжать, отбывать
leave	оставлять, покидать, уходить, уезжать, бросать, удаляться
exit	выходить, уходить, умереть
get out	выходить, выбираться, вылезать, убираться, уходить, выкарабкиваться
get away	избежать, уходить, ускользнуть, удирать, выбираться, отрываться
walk away	уходить, унести, обходить стороной, уводить, украсть
walk off	уходить, унести, одержать легкую победу, украсть, уводить
go off	уходить, сходить, выстрелить, уезжать, сбежать, проходить
withdraw	изымать, отзываться, уходить, забирать, увести, удаляться
take away	увести, забрать, забирать, отнимать, унести, уходить
come away	уходить, отламываться, отходить
retire	уходить, удаляться, уходить в отставку, увольняться, отступать, уединяться
retreat	отступать, отходить, уходить, удаляться, отбросить назад
walk out	уходить, демонстративно покинуть, выходить
be off	уходить
buzz off	уходить, удаляться, смыться, улизнуть
scram	уходить, удирать
be away	отсутствовать, уходить, не быть дома
shove off	убираться, уходить, отталкиваться
toddle off	уходить
toddle	ковылять, учиться ходить, прогуливаться, бродить, уходить
draw away	отвлекать, удаляться, уводить, уходить, отрываться от противника
split	раскалывать, разбивать, расцепить, раскалываться, расцепляться, уходи

Статистический перевод по словам — Word-based SMT

Model 1: мешок слов

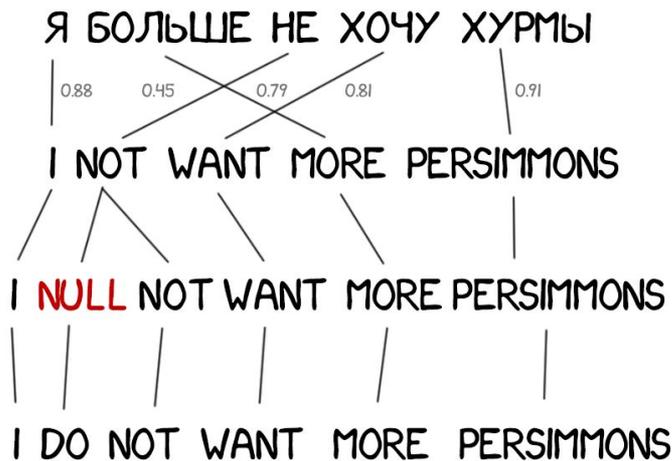
Классический подход — делим всё на слова и считаем статистику. Никакого учёта порядка или перестановок. Из хитростей Model 1 умела разве что переводить одно слово в несколько. Der Staubsauger (пылесос) легко превращался в Vacuum Cleaner, но обратно уже как повезет.

Я БОЛЬШЕ НЕ ХОЧУ ХУРМЫ
| 0.88 / 0.45 / 0.79 / 0.81 | 0.91
I MORE NOT WANT PERSIMMONS

Model 3: добавление отсутствующих слов

Часто при переводе появляются новые слова, которых не было в оригинальном тексте. В немецком языке внезапно вылезают артикли, в английском вставляют глагол *do* где не попадя. «Я не хочу хурмы» → «I **do** not want persimmons. Чтобы решить эту проблему в Model 3 добавили два промежуточных шага:

1. Вставка маркеров (NULL-слов) на те места, где машина подозревает необходимость нового слова
2. Подбор нужного артикля, частицы или глагола под каждый маркер



Model 4: перестановки слов

Model 2 хоть учитывала порядок слов в предложении, но ничего не знала про перестановки слов между собой. Часто при переводе надо, например, поменять существительное и прилагательное местами. Тут сколько ни запоминай их порядок по всему предложению — лучше не станет. Потому в Model 4 стали учитывать еще и так называемый «относительный порядок». Если при переводе два слова постоянно менялись друг с другом — модель это запоминала.

Статистический перевод по фразам — Phrase-based SMT

Для обучения он разбивал текст не только на слова, но и на целые фразы. Точнее N-граммы или фраземы — пересекающиеся наборы из N слов подряд. Машина научилась переводить устойчивые сочетания слов, что заметно улучшило точность.

МОЖЕТ ХВАТИТ ПРИМЕРОВ С ХУРМОЙ

УНИГРАММЫ:

1. МОЖЕТ
2. ХВАТИТ
3. ПРИМЕРОВ
4. С
5. ХУРМОЙ

МОЖЕТ ХВАТИТ ПРИМЕРОВ С ХУРМОЙ

БИГРАММЫ:

1. МОЖЕТ ХВАТИТ
2. ХВАТИТ ПРИМЕРОВ
3. ПРИМЕРОВ С
4. С ХУРМОЙ

МОЖЕТ ХВАТИТ ПРИМЕРОВ С ХУРМОЙ

ТРИГРАММЫ:

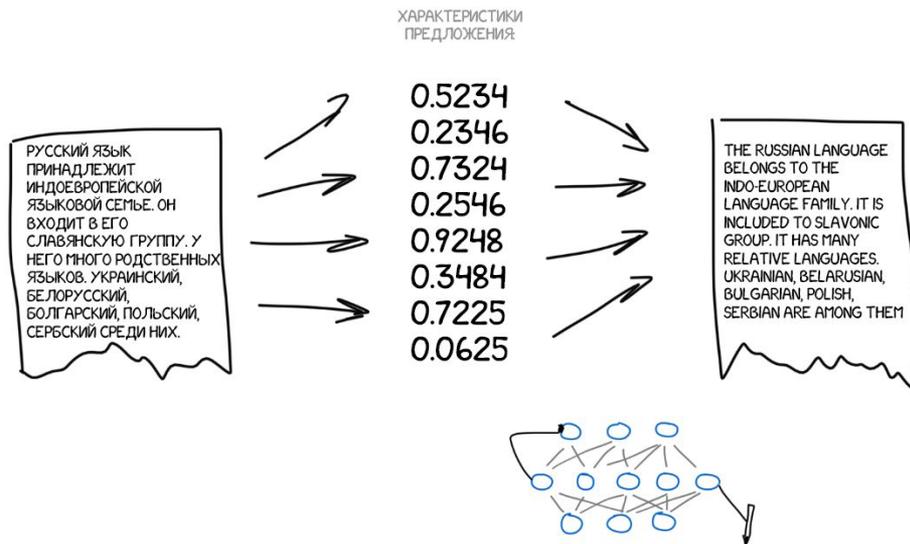
1. МОЖЕТ ХВАТИТ ПРИМЕРОВ
2. ХВАТИТ ПРИМЕРОВ С
3. ПРИМЕРОВ С ХУРМОЙ

Нейронный машинный перевод — Neural Machine Translation (NMT)

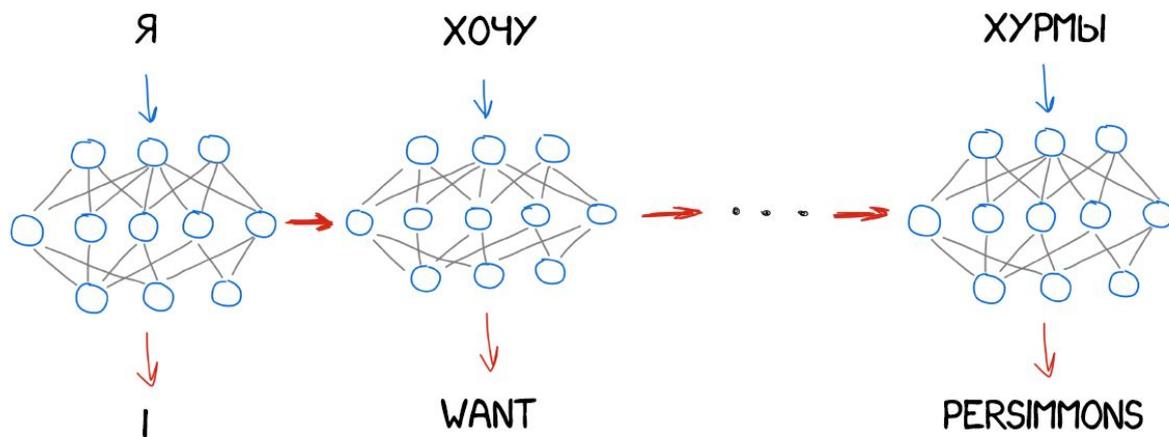
Помните приложение Prisma, которое обрабатывало фото в стиле известного художника? Там не было особой магии — нейросеть обучили распознавать картины художника, а потом «оторвали» последние слои, где она принимает решение. Получившиеся наброски, по сути промежуточное представление сети, и было той самой стилизованной картинкой.



Первая нейросеть умеет только кодировать предложение в набор цифр-характеристик, а вторая только декодировать их обратно в текст. Обе понятия не имеют друг о друге, каждая знает только свой язык.

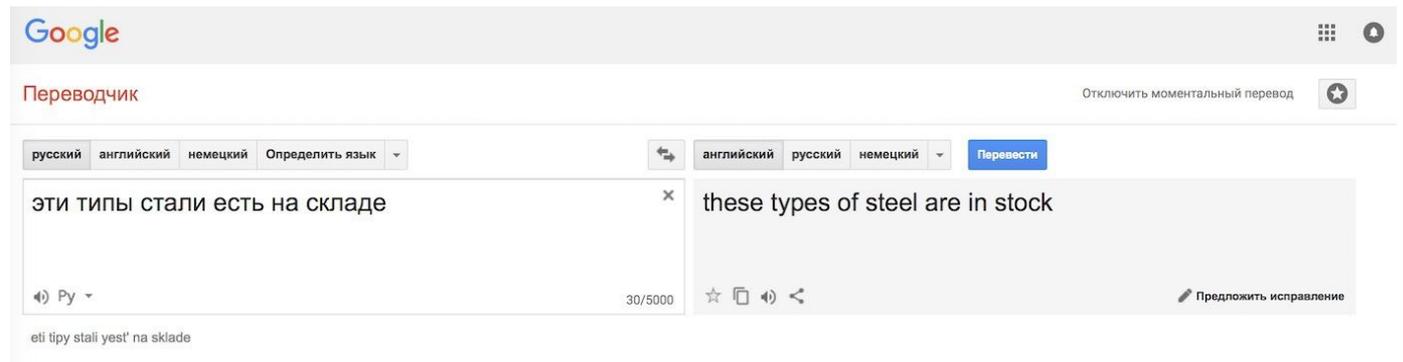


RNN сейчас применяют в: распознавание речи в Siri , подсказки слов на клавиатуре , генерация музыки и даже чатботы.



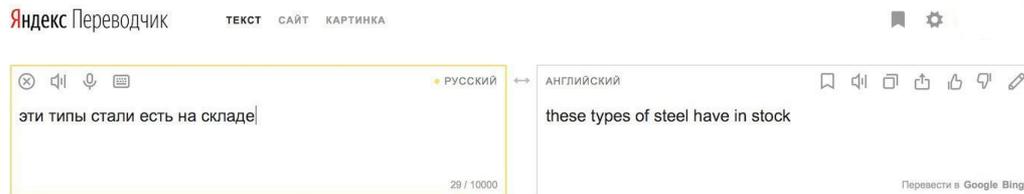
Google Translate (2016)

В 2016 году Google включил нейронный перевод девяти языков между собой, в 2017 был добавлен и русский. Google разработал собственную систему под нехитрым названием Google Neural Machine Translation (GNMT), состоявшую аж из 8-слойного RNN на входе и такого же на выходе и системы согласования контекста под названием Attention Model.

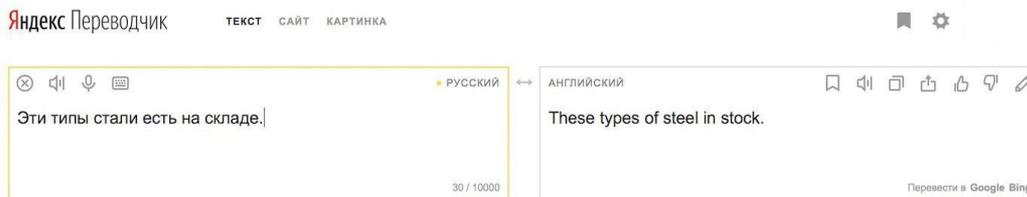


Яндекс Переводчик (2017)

Яндекс запустил свой нейросетевой перевод в 2017 году. Главным отличием они заявили гибридность. Переводчик Яндекса переводит предложение сразу двумя методами — статистическим и нейросетевым, а потом с помощью их любимого алгоритма CatBoost находит наиболее подходящий.



ДОБАВЛЯЕМ ТОЧКУ В КОНЦЕ ПРЕДЛОЖЕНИЯ, ЯНДЕКС ВКЛЮЧАЕТ НЕЙРОСЕТИ И НАЧИНАЕТ ПЕРЕВОДИТЬ ЛУЧШЕ
ТАКОЕ ВОТ МАШИННОЕ ОБУЧЕНИЕ



Заключение и будущее

Всех по-прежнему будоражит идея «Вавилонской Рыбки» — синхронного перевода речи на лету. Google делала шаг в этом направлении, когда анонсировала Pixel Buds, но на поверку всё оказалось плохо. Синхронный перевод на лету отличается от обычного, ведь нужно знать места, когда начать переводить, а когда сидеть и слушать. Подходов к решению этой задачи найти не удалось.