

アダプテーションによる 製薬専門NMTの研究開発

情報通信研究機構

内山将夫

JTF関西セミナー

2020年6月23日

概要

- NMTの概要とみんなの自動翻訳@TexTra
- 製薬専門NMT

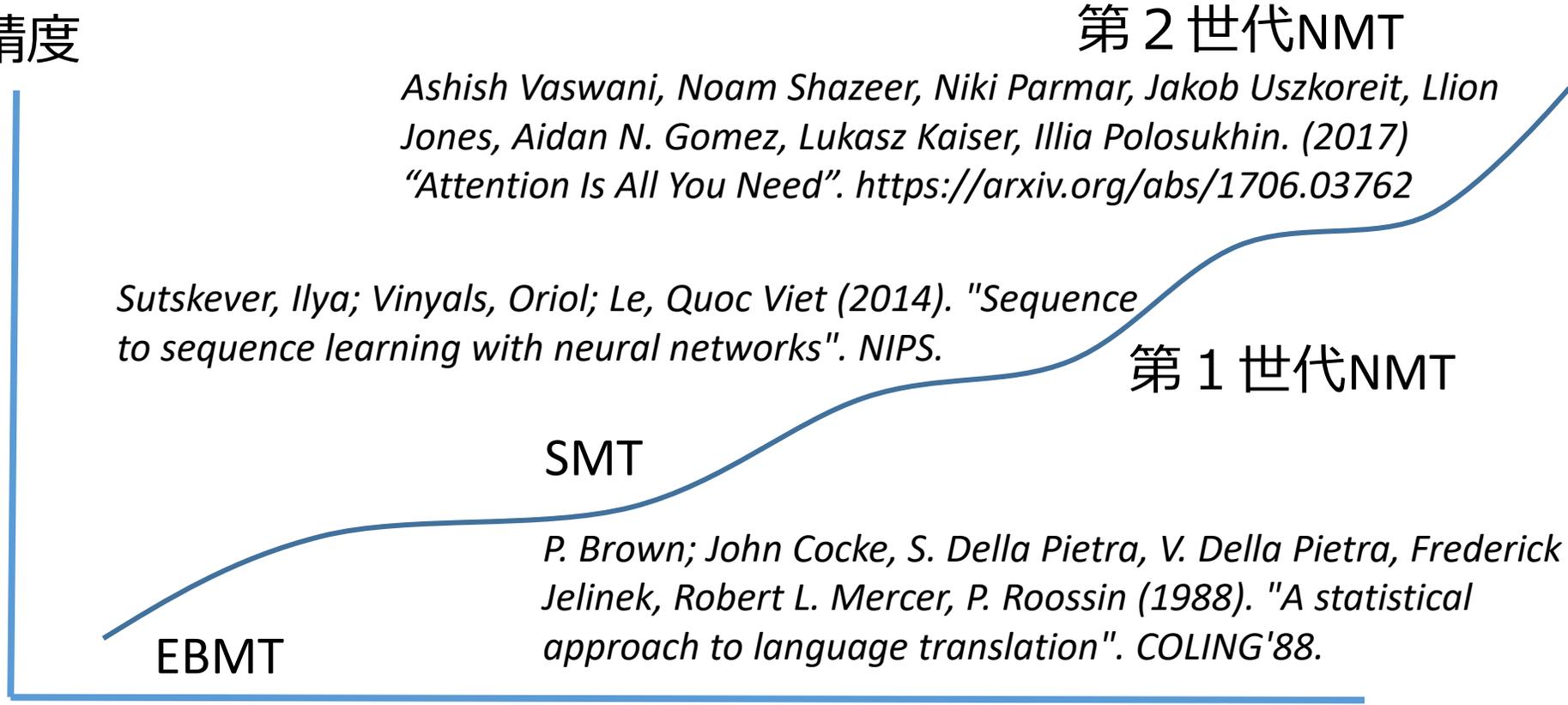
コーパスベースMTの構成要素

- 自動翻訳アルゴリズム
1980年代からの発展⇒NMT
- 学習データとしての対訳コーパス
異なる言語で同一意味の文章の組からなるデータベース
- MTの評価手法
人間がどのように感じるかの観点からのMTの良さ
- MTの使い方
MTの普及により、MTの適切な使い方の周知が必要



コーパスベースMTアルゴリズムの進展

翻訳精度



Makoto Nagao (1984). "A framework of a mechanical translation between Japanese and English by analogy principle". In A. Elithorn and R. Banerji. Artificial and Human Intelligence. Elsevier Science Publishers

年代

第2世代NMT（汎用NT） @みんなの自動翻訳



The screenshot shows a web browser window with the URL <https://mt-auto-minhon-mlt.ucrj.jgn-x.jp/content/demo/>. The page title is "みんなの自動翻訳@TexTra®". The interface includes a navigation menu with "翻訳データ", "自動翻訳", and "ツール". A "ヘルプ" button is visible in the top right.

The main content area is titled "自動翻訳". It features a language selection interface with buttons for "英語", "日本語", and "汎用NT 【英語 - 日本語】 1". A red box highlights the "汎用NT" button. A green "翻訳" button is also present.

The source text (English) is displayed in a box on the left:

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder-decoder configuration. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train.

The translated text (Japanese) is displayed in a box on the right:

支配的なシーケンス変換モデルは、エンコーダ/デコーダ構成における複雑なリカレントまたは畳み込みニューラルネットワークに基づいている。また、最適なモデルは、アテンション機構を介してエンコーダとデコーダを接続する。本稿では、注意機構のみに基づき、リカレントと畳み込みを完全に省く、新しい簡単なネットワークアーキテクチャTransformerを提案する。2つの機械翻訳タスクに関する実験では、これらのモデルは、より並列化可能であり、トレーニングに要する時間が大幅に短縮される一方で、品質が優れていることが示されている。

At the bottom, there are checkboxes for "文章を結合する" and "文章を区切る", and buttons for "← 訳文コピー" and "× リセット".

汎用NT @みんなの自動翻訳

原文英語	汎用NT (第2世代)
<p>The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder-decoder configuration.</p> <p>The best performing models also connect the encoder and decoder through an attention mechanism.</p> <p>We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely.</p> <p>Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train.</p>	<p>支配的なシーケンス変換モデルは、エンコーダ/デコーダ構成における複雑なリカレントまたは畳み込みニューラルネットワークに基づいている。</p> <p>また、最適なモデルは、アテンション機構を介してエンコーダとデコーダを接続する。</p> <p>本稿では、注意機構のみに基づき、リカレントと畳み込みを完全に省く、新しい簡単なネットワークアーキテクチャTransformerを提案する。</p> <p>2つの機械翻訳タスクに関する実験では、これらのモデルは、より並列化可能であり、トレーニングに要する時間が大幅に短縮される一方で、品質が優れていることが示されている。</p>

*English sentences are part of Abstract from Ashish Vaswani, et al. (2017) "Attention Is All You Need".
<https://arxiv.org/abs/1706.03762>*

汎用NT @みんなの自動翻訳

原文日本語	汎用NT（第2世代）
<p>2016年のニューラル機械翻訳（NMT）の実用化は、翻訳業界に衝撃を与え、ポケトークのような自動翻訳端末の市場拡大につながるなど、社会に大きなインパクトを与えた。</p> <p>ただし、翻訳技術や自然言語処理技術（NLP）分野では、その後も革命級のブレークスルーが相次いでいる。</p> <p>翻訳を含む言語系の人工知能（AI）が従来の常識を次々と塗り替え、ありえないペースで発展している。</p>	<p>The practical application of Neural Machine Translation (NMT) in 2016 gave a shock to the translation industry and had a large impact on society, such as the expansion of the market for automatic translation terminals such as PokeTalk.</p> <p>However, breakthroughs in translation and natural language processing (NLP) have continued.</p> <p>Artificial intelligence (AI), a linguistic system that includes translation, is evolving at an incredible pace, breaking conventional wisdom.</p>

原文は次の一部：野澤 哲生「AI翻訳が人間超え、言葉の壁崩壊へ」日経エレクトロニクス、2019/08/20
<https://tech.nikkeibp.co.jp/atcl/nxt/mag/ne/18/00046/00002/>

自動評価尺度 BLEU

- M T 訳と参照訳の類似度の尺度
- 1 0 0 0 文～5 0 0 0 文程度の文章で計算
- 一般的に BLEU 3 0 ～4 0 以上であれば高精度

Papineni, Kishore and Roukos, Salim and Ward, Todd and Zhu, Wei-Jing. (2002)
BLEU: a Method for Automatic Evaluation of Machine Translation. ACL

みんなの自動翻訳でのBLEU計算法

- 自動翻訳 > 自動翻訳評価 > 「+新規登録」



The screenshot shows a web browser window with the URL <https://mt-auto-minhon-mlt.ucrjgn-x.jp/content/eval/edit/index.html>. The page title is "みんなの自動翻訳@TexTra®". The main heading is "自動翻訳評価 登録". The form includes a name input field with the placeholder "名前を入力してください。", a language direction selector set to "日本語 ↔ 英語", and an automatic translation mode selector set to "標準". Below these are several checkboxes for different evaluation methods, all currently unchecked:

- 汎用NT【日本語 - 英語】 ⓘ
- 汎用NMT【日本語 - 英語】 ⓘ
- 特許NT【日本語 - 英語】 ⓘ
- 特許NMT【日本語 - 英語】 ⓘ
- 特許請求項NMT【日本語 - 英語】 ⓘ
- 対話NMT【日本語 - 英語】 ⓘ
- ミニ対話NMT【日本語 - 英語】 ⓘ

第2世代製薬NMTは高精度

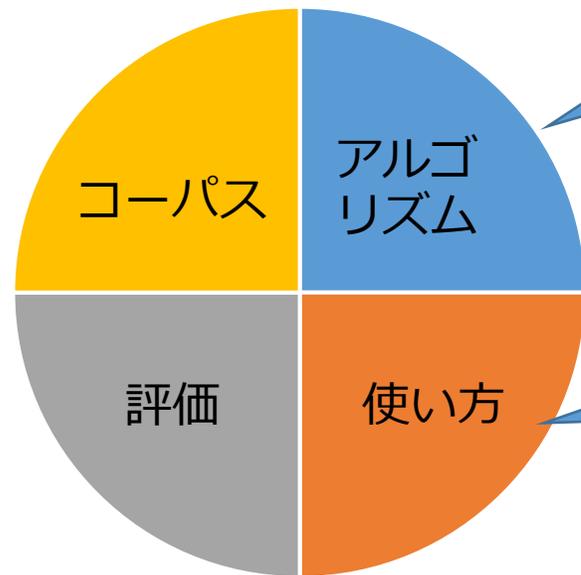
	第1世代	第2世代
汎用	23.8	30.3
製薬	39.9	49.4

第1世代 << 第2世代

汎用 << 製薬

自動翻訳尺度BLEUの比較

NMTにより新しくなった点



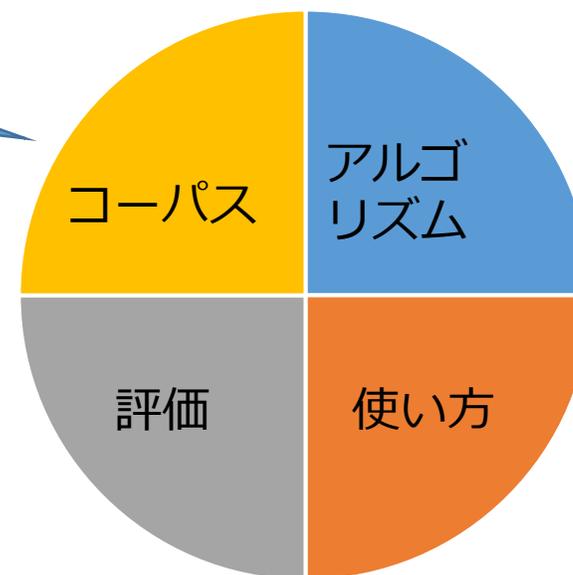
ニューラルネットワークによる高精度アルゴリズム

高精度自動翻訳の社会的普及

NMTでも以前と同じ点

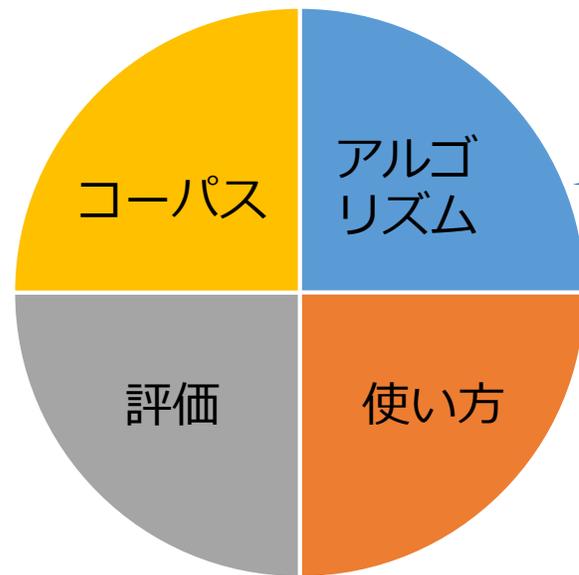
大規模対訳
コーパスが必要

人間による評価
が必要



NMTアルゴリズムの概要

- 基本訓練
- アダプテーション



ニューラルネットワークによる高精度
アルゴリズム

NMTの訓練とアダプテーション概要

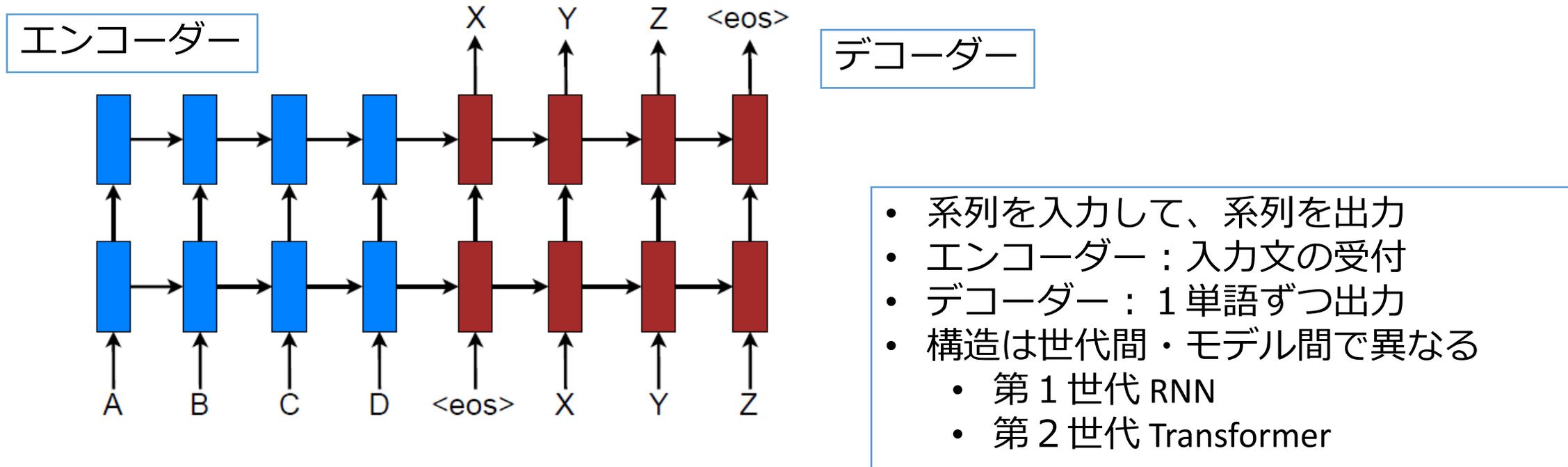
【訓練】 1文ずつNMTモデルのパラメタを調整する

- 大雑把には：
 - 入力文を翻訳
 - 参照訳文と比較
 - 翻訳文と参照訳文の違いに応じてNMTのパラメタを更新
 - 以上を大規模に繰り返す（数億回になることもある）

【アダプテーション】

- 訓練済みNMTモデルに、上記訓練を特定分野データで追加
- 訓練済みモデルをベースにするので、比較的少量データで高精度

NMTの基本構造 (第1・2世代ほぼ共通)



- 系列を入力して、系列を出力
- エンコーダー：入力文の受付
- デコーダー：1単語ずつ出力
- 構造は世代間・モデル間で異なる
 - 第1世代 RNN
 - 第2世代 Transformer

Figure 1: **Neural machine translation** – a stacking recurrent architecture for translating a source sequence $A B C D$ into a target sequence $X Y Z$. Here, $\langle eos \rangle$ marks the end of a sentence.

Thang Luong; Hieu Pham; Christopher D. Manning. (2015)
Effective Approaches to Attention-based Neural Machine Translation. EMNLP

第1世代と第2世代の精度比較

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

翻訳の比較的難しい英独で BLEU 4ポイント程度の大きな改善

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin. (2017) "Attention Is All You Need". <https://arxiv.org/abs/1706.03762>

アダプテーション@みんなの自動翻訳



自分好みの自動翻訳イン: × + ▾

← → ↻ 🏠 🔒 <https://mt-auto-minhon-mlt.ucr.jgn-x.jp/content/introduction/adapt.html> 📖 ☆ ⚙️ 🖨️ 🔗 ⋮

🍃 原文

Download the All download option.

🗨️ 汎用NMT

downloadオプションをダウンロードする。(少し残念な訳)

🗨️ アダプテーションMT (2千対訳文)

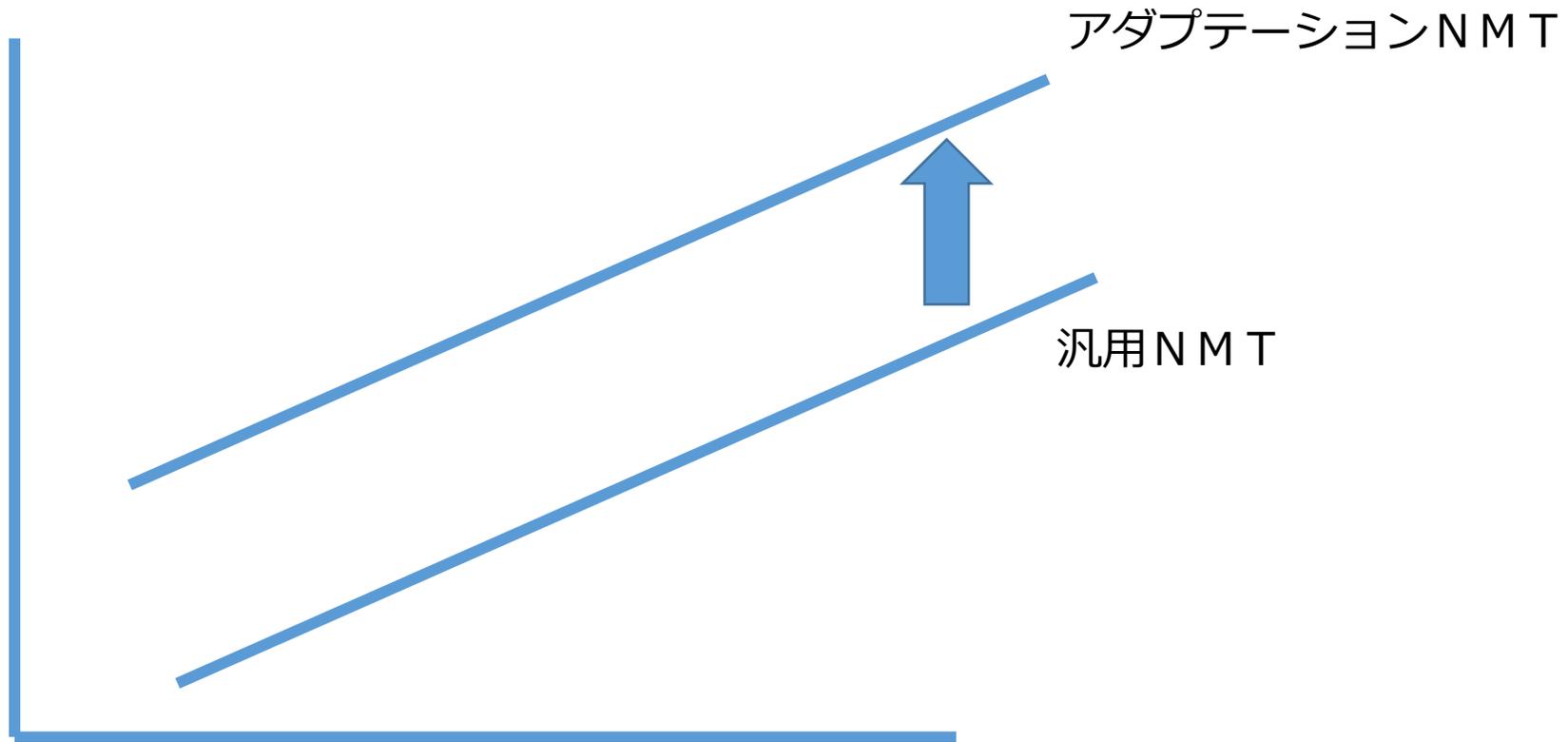
ダウンロード・オプションをすべてダウンロード。(誤訳)

🗨️ アダプテーションMT (1万対訳文)

ダウンロードオプション「All」をダウンロードする。(完璧！👍)

汎用とアダプテーションNMTの関係

翻訳精度



対訳量

製薬NMT : R&D Head Club + NICT

- RDHCの8社から提供された320万文対以上の日英対訳データ
- 一昨年度のNICTとアストラゼネカの協業を発展
- 汎用NTを製薬対訳でアダプテーション
- 事業会社によりサービス提供が開始済み

第2世代製薬NMTは高精度（再掲）

	第1世代	第2世代
汎用	23.8	30.3
製薬	39.9	49.4

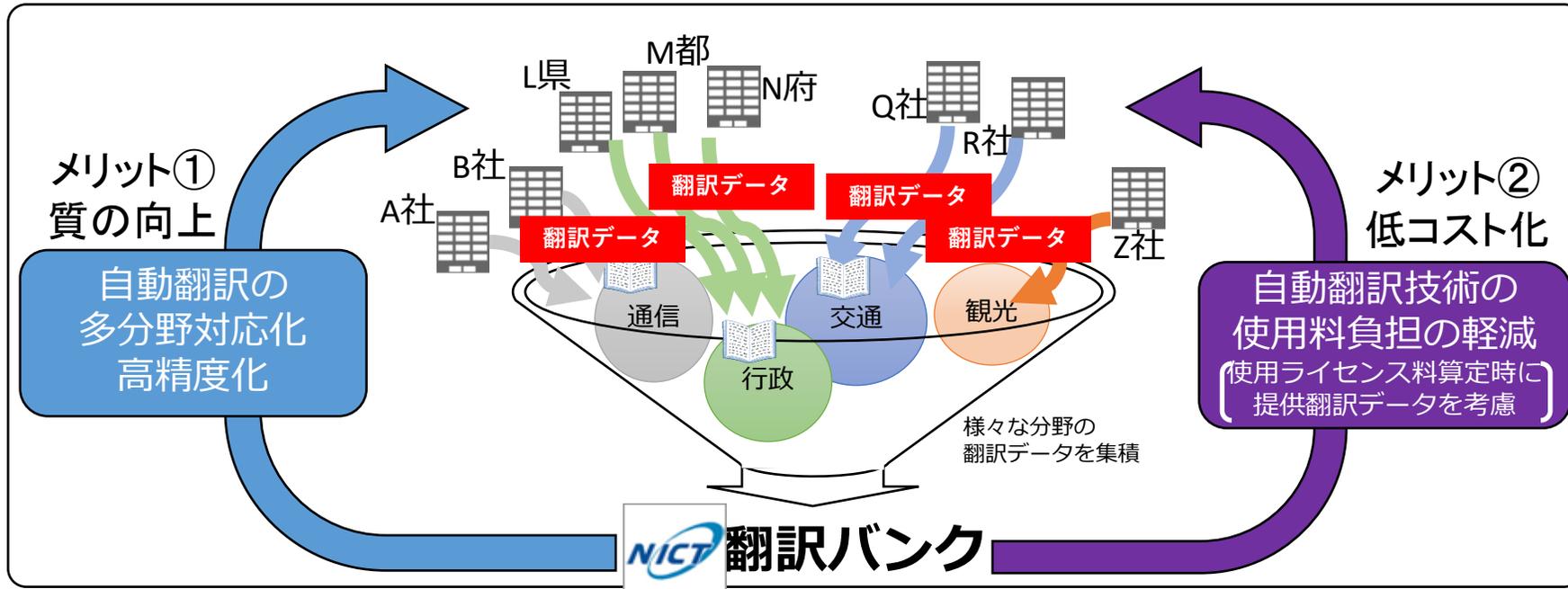
第1世代 << 第2世代

汎用 << 製薬

自動翻訳尺度BLEUの比較

翻訳バンクで高品質対訳 アダプテーションで高精度NMTを実現

NICTの自動翻訳技術の使用ライセンス料の算定の際に、提供が見込まれる翻訳データを勘案して負担を軽減する仕組みを導入



覚えておいていただきたいこと

- 自動翻訳エンジンは進化します
- アルゴリズムの進化
- 対訳データの進化
- 自動翻訳精度の向上に頭打ち感はありません
- 自動翻訳をみんなで育てましょう！