

Recent Advances In Document-level Neural Machine Translation

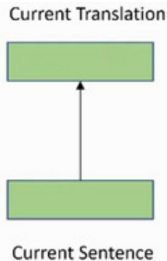
Lorenzo Lupo

Supervisors: Laurent Besacier, Marco Dinarelli

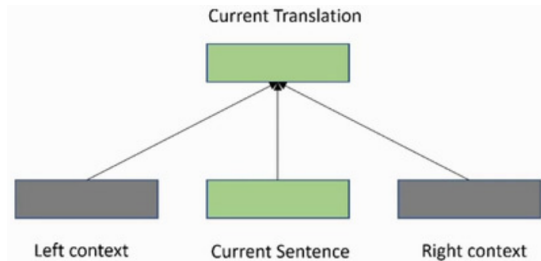
July 10, 2020

What is Document-level Machine Translation

Sentence-level MT

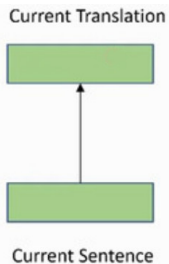


Document-level MT

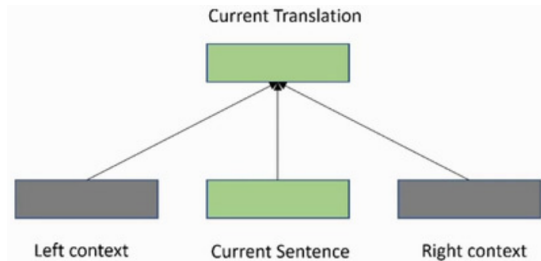


Document-level MT \leftrightarrow Context-aware MT

Context-agnostic MT



Context-aware MT



Why Document-level NMT ?

Why Document-level NMT ?

- ▶ Some recent results suggest that neural machine translation (NMT) "approaches the accuracy achieved by average bilingual human translators [on some test sets] [Wu et al., 2016]"

Why Document-level NMT ?

- ▶ Some recent results suggest that neural machine translation (NMT) "approaches the accuracy achieved by average bilingual human translators [on some test sets] [Wu et al., 2016]"
- ▶ "In a pairwise ranking experiment, human raters assessing **adequacy** and **fluency** show a stronger preference for human over machine translation when evaluating documents as compared to isolated sentences." [Lubli et al., 2018]

Sentence-level NMT is inconsistent

B: How are you today?

Sentence-level NMT is inconsistent

B: How are *you* today?

SENTENCE-LEVEL TRANSLATION

B: Comment *vas-tu* aujourd'hui ?

Sentence-level NMT is inconsistent

A: Good Morning, Mr. President.

B: How are you today?

SENTENCE-LEVEL TRANSLATION

B: Comment vas-tu aujourd'hui ?

Sentence-level NMT is inconsistent

A: Good Morning, Mr. President.

B: How are you today?

SENTENCE-LEVEL TRANSLATION

B: Comment vas-tu aujourd'hui ?

DOCUMENT-LEVEL TRANSLATION

B: Comment allez-vous aujourd'hui ?

How frequent are inconsistencies ?

[[Voita et al., 2019b](#)] undertake a human study on context agnostic translation :

- 2000 pairs of consecutive English sentences ($S1 + S2$) from OpenSubtitles2018
- translate to Russian with Transformer model [[Vaswani et al., 2017](#)]

How frequent are inconsistencies ?

[[Voita et al., 2019b](#)] undertake a human study on context agnostic translation :

- 2000 pairs of consecutive English sentences (S1 + S2) from OpenSubtitles2018
- translate to Russian with Transformer model [[Vaswani et al., 2017](#)]

all	one/both bad	both good	
		bad pair	good pair
2000	211	140	1649
100%	11%	7%	82%

Which kind of inconsistencies?

type of phenomena	frequency
deixis	37%
ellipsis	29%
lexical cohesion	14%
ambiguity	9%
anaphora	6%
other	5%

Figure: Types of phenomena causing inconsistencies between English-Russian context-agnostic translations of consecutive sentences when placed in the context of each other.

Objectives

Objectives

- ▶ **Design translation models and learning techniques** that solve inconsistencies by taking context into account;

Objectives

- **Design translation models and learning techniques** that solve inconsistencies by taking context into account;
- **Evaluate such models** in a proper way;

Plan

1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

Plan

1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

- ▶ Classical metrics such as BLUE and METEOR are inadequate in evaluating document-level neural MT (DNMT) because they evaluate **average translation quality** by matching n-grams at **sentence-level**. Thus:

- Classical metrics such as BLUE and METEOR are inadequate in evaluating document-level neural MT (DNMT) because they evaluate **average translation quality** by matching n-grams at **sentence-level**. Thus:
 - they are unable to capture document-wide phenomena like **coherence** and **cohesion** [Wong and Kit, 2012].

- Classical metrics such as BLUE and METEOR are inadequate in evaluating document-level neural MT (DNMT) because they evaluate **average translation quality** by matching n-grams at **sentence-level**. Thus:
 - they are unable to capture document-wide phenomena like **coherence** and **cohesion** [Wong and Kit, 2012].
 - they are not able to measure improvements over translations of discourse **phenomena that affect few words but heavily influence fluency and correctness** of the translation [Miller et al., 2018]. E.g. pronominal anaphora.

- Classical metrics such as BLUE and METEOR are inadequate in evaluating document-level neural MT (DNMT) because they evaluate **average translation quality** by matching n-grams at **sentence-level**. Thus:
 - they are unable to capture document-wide phenomena like **coherence** and **cohesion** [Wong and Kit, 2012].
 - they are not able to measure improvements over translations of discourse phenomena that affect few words but heavily influence fluency and correctness of the translation [Miller et al., 2018]. E.g. pronominal anaphora.
- Instead, we should evaluate DNMT with metrics that can **capture inter-sentential discourse phenomena**.

Desiderata of Document-level MT Evaluation

The evaluation of document-level MT should:

¹in the remainder of this presentation, we refer to inter-sentential context simply as context.

Desiderata of Document-level MT Evaluation

The evaluation of document-level MT should:

- Be undertaken by providing inter-sentential context¹.

¹in the remainder of this presentation, we refer to inter-sentential context simply as context.

Desiderata of Document-level MT Evaluation

The evaluation of document-level MT should:

- Be undertaken by providing inter-sentential context¹.
- Consider both translation quality (BLEU) and correct translation of discourse-phenomena.

¹in the remainder of this presentation, we refer to inter-sentential context simply as context.

Desiderata of Document-level MT Evaluation

The evaluation of document-level MT should:

- Be undertaken by providing inter-sentential context¹.
- Consider both translation quality (BLEU) and correct translation of discourse-phenomena.
- Focus on context-dependent discourse phenomena.

¹in the remainder of this presentation, we refer to inter-sentential context simply as context.

Desiderata of Document-level MT Evaluation

The evaluation of document-level MT should:

- Be undertaken by providing inter-sentential context¹.
- Consider both translation quality (BLEU) and correct translation of discourse-phenomena.
- Focus on context-dependent discourse phenomena.
 - E.g., pronominal anaphora cases in which the antecedent is in a previous sentence (context-dependent), instead of being in the same sentence (context-independent).

¹in the remainder of this presentation, we refer to inter-sentential context simply as context.

Desiderata of Document-level MT Evaluation

The evaluation of document-level MT should:

- Be undertaken by providing inter-sentential context¹.
- Consider both translation quality (BLEU) and correct translation of discourse-phenomena.
- Focus on context-dependent discourse phenomena.
 - E.g., pronominal anaphora cases in which the antecedent is in a previous sentence (context-dependent), instead of being in the same sentence (context-independent).
- Focus on hard cases.

¹in the remainder of this presentation, we refer to inter-sentential context simply as context.

Desiderata of Document-level MT Evaluation

The evaluation of document-level MT should:

- Be undertaken by providing inter-sentential context¹.
- Consider both translation quality (BLEU) and correct translation of discourse-phenomena.
- Focus on context-dependent discourse phenomena.
 - E.g., pronominal anaphora cases in which the antecedent is in a previous sentence (context-dependent), instead of being in the same sentence (context-independent).
- Focus on hard cases.
 - E.g., when translating English to French, **he** is easy whereas **it** is hard to translate because ambiguous.

¹in the remainder of this presentation, we refer to inter-sentential context simply as context.

Plan

1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

Automatic metrics

In the literature, we can find some automatic metrics for machine translation that specialize on discourse phenomena:

Automatic metrics

In the literature, we can find some automatic metrics for machine translation that specialize on discourse phenomena:

- Metrics for **pronoun translation**, usually involved in phenomena of coreference resolution (deixis, grammatical cohesion). E.g.,

Automatic metrics

In the literature, we can find some automatic metrics for machine translation that specialize on discourse phenomena:

- Metrics for **pronoun translation**, usually involved in phenomena of coreference resolution (deixis, grammatical cohesion). E.g.,
 - **Accuracy of Pronoun Translation** [Miculicich Werlen and Popescu-Belis, 2017]
and **Pronoun Pair-wise Ranking** [Jwalapuram et al., 2019]

Automatic metrics

In the literature, we can find some automatic metrics for machine translation that specialize on discourse phenomena:

- Metrics for **pronoun translation**, usually involved in phenomena of coreference resolution (deixis, grammatical cohesion). E.g.,
 - **Accuracy of Pronoun Translation** [[Miculicich Werlen and Popescu-Belis, 2017](#)]
and **Pronoun Pair-wise Ranking** [[Jwalapuram et al., 2019](#)]
- Metrics for **lexical cohesion**. E.g.,

Automatic metrics

In the literature, we can find some automatic metrics for machine translation that specialize on discourse phenomena:

- Metrics for **pronoun translation**, usually involved in phenomena of coreference resolution (deixis, grammatical cohesion). E.g.,
 - **Accuracy of Pronoun Translation** [Miculicich Werlen and Popescu-Belis, 2017] and **Pronoun Pair-wise Ranking** [Jwalapuram et al., 2019]
- Metrics for **lexical cohesion**. E.g.,
 - [Hajlaoui and Popescu-Belis, 2013] proposed new automatic and semi-automatic metrics for discourse connectives, referred to as **Accuracy of Connective Translation**.

Automatic metrics

In the literature, we can find some automatic metrics for machine translation that specialize on discourse phenomena:

- Metrics for **pronoun translation**, usually involved in phenomena of coreference resolution (deixis, grammatical cohesion). E.g.,
 - **Accuracy of Pronoun Translation** [Miculicich Werlen and Popescu-Belis, 2017] and **Pronoun Pair-wise Ranking** [Jwalapuram et al., 2019]
- Metrics for **lexical cohesion**. E.g.,
 - [Hajlaoui and Popescu-Belis, 2013] proposed new automatic and semi-automatic metrics for discourse connectives, referred to as **Accuracy of Connective Translation**.
 - [Wong and Kit, 2012] measure the abundance of lexical cohesion characterized by semantically related words with **Lexical Cohesion Devices**.

Lexical Cohesion Devices [[Wong and Kit, 2012](#)]

Lexical Cohesion Devices [Wong and Kit, 2012]

1. A **stemming algorithm** [Porter, 1980] is used to identify word stems for each content word .

Lexical Cohesion Devices [Wong and Kit, 2012]

1. A **stemming algorithm** [Porter, 1980] is used to identify word stems for each content word .
 - Words with the same stem are defined and counted as **Repetitions**.

Lexical Cohesion Devices [Wong and Kit, 2012]

1. A **stemming algorithm** [Porter, 1980] is used to identify word stems for each content word .
 - Words with the same stem are defined and counted as **Repetitions**.
2. **WordNet** [Fellbaum, 1998] is used to cluster synonyms and superordinates into semantic groups.

Lexical Cohesion Devices [Wong and Kit, 2012]

1. A **stemming algorithm** [Porter, 1980] is used to identify word stems for each content word .
 - Words with the same stem are defined and counted as **Repetitions**.
2. **WordNet** [Fellbaum, 1998] is used to cluster synonyms and superordinates into semantic groups.
 - Words belonging to the same semantic group or close semantic groups (near-synonyms) are defined and counted as **Lexical Cohesion Devices (LCD)**.

Lexical Cohesion Devices [Wong and Kit, 2012]

1. A **stemming algorithm** [Porter, 1980] is used to identify word stems for each content word .
 - Words with the same stem are defined and counted as **Repetitions**.
2. **WordNet** [Fellbaum, 1998] is used to cluster synonyms and superordinates into semantic groups.
 - Words belonging to the same semantic group or close semantic groups (near-synonyms) are defined and counted as **Lexical Cohesion Devices** (LCD).
3. A **hybrid metric** can then be defined as weighted average of:

Lexical Cohesion Devices [Wong and Kit, 2012]

1. A **stemming algorithm** [Porter, 1980] is used to identify word stems for each content word .
 - Words with the same stem are defined and counted as **Repetitions**.
2. **WordNet** [Fellbaum, 1998] is used to cluster synonyms and superordinates into semantic groups.
 - Words belonging to the same semantic group or close semantic groups (near-synonyms) are defined and counted as **Lexical Cohesion Devices** (LCD).
3. A **hybrid metric** can then be defined as weighted average of:
 - a **classic sentence-level metric**, e.g. BLEU, METEOR, TER.

Lexical Cohesion Devices [Wong and Kit, 2012]

1. A **stemming algorithm** [Porter, 1980] is used to identify word stems for each content word .
 - Words with the same stem are defined and counted as **Repetitions**.
2. **WordNet** [Fellbaum, 1998] is used to cluster synonyms and superordinates into semantic groups.
 - Words belonging to the same semantic group or close semantic groups (near-synonyms) are defined and counted as **Lexical Cohesion Devices** (LCD).
3. A **hybrid metric** can then be defined as weighted average of:
 - a **classic sentence-level metric**, e.g. BLEU, METEOR, TER.
 - a **lexical cohesion metric**, e.g. *Repetitions/content words* or *LCD/content words*.

Lexical Cohesion Devices [Wong and Kit, 2012]

1. A **stemming algorithm** [Porter, 1980] is used to identify word stems for each content word .
 - Words with the same stem are defined and counted as **Repetitions**.
 2. **WordNet** [Fellbaum, 1998] is used to cluster synonyms and superordinates into semantic groups.
 - Words belonging to the same semantic group or close semantic groups (near-synonyms) are defined and counted as **Lexical Cohesion Devices** (LCD).
 3. A **hybrid metric** can then be defined as weighted average of:
 - a **classic sentence-level metric**, e.g. BLEU, METEOR, TER.
 - a **lexical cohesion metric**, e.g. *Repetitions/content words* or *LCD/content words*.
- *Compatible languages*: all languages with stemmers and WordNets available.

Plan

1. Evaluation

1.2 Automatic metrics

1.3 Test Suites

1.4 Remarks

2. Approaches to DNMT

2.5 Concatenation Approaches

2.6 Separate Encoding Approaches

2.7 Cache Approaches

2.8 Exploiting Document-level Monolingual Corpora

2.9 Others

2.10 Remarks and conclusions

Test Suites

In the literature we can distinguish three kinds of test suites:

Test Suites

In the literature we can distinguish three kinds of test suites:

- **Manual test suites** consist in the manual evaluation of a number of test cases (e.g. discourse phenomena) based on a given machine translation task.

Test Suites

In the literature we can distinguish three kinds of test suites:

- **Manual test suites** consist in the manual evaluation of a number of test cases (e.g. discourse phenomena) based on a given machine translation task.
 - For instance, WMT2019 not only provided ratings for each system output but also some manual evaluation and analysis of the outputs, like the English → Czech test suite on coherence by [Rysov et al., 2019].

Test Suites

In the literature we can distinguish three kinds of test suites:

- ▶ **Manual test suites** consist in the manual evaluation of a number of test cases (e.g. discourse phenomena) based on a given machine translation task.
 - ▶ For instance, WMT2019 not only provided ratings for each system output but also some manual evaluation and analysis of the outputs, like the English → Czech test suite on coherence by [Rysov et al., 2019].
- ▶ **Specialized test sets** are like normal MT test sets but consist of consecutive sentences that are more densely populated with specific discourse phenomena. Translations are evaluated on such tests sets by means of average quality metrics like BLEU.

Test Suites

In the literature we can distinguish three kinds of test suites:

- ▶ **Manual test suites** consist in the manual evaluation of a number of test cases (e.g. discourse phenomena) based on a given machine translation task.
 - ▶ For instance, WMT2019 not only provided ratings for each system output but also some manual evaluation and analysis of the outputs, like the English → Czech test suite on coherence by [Rysov et al., 2019].
- ▶ **Specialized test sets** are like normal MT test sets but consist of consecutive sentences that are more densely populated with specific discourse phenomena. Translations are evaluated on such tests sets by means of average quality metrics like BLEU.
 - ▶ E.g. [Voita et al., 2018] build a specialized English → Russian test set by retrieving from OpenSubtitles2016 all the sentences containing pronouns that are coreferent to an expression in the previous sentence.

Test Suites

In the literature we can distinguish three kinds of test suites:

- ▶ **Manual test suites** consist in the manual evaluation of a number of test cases (e.g. discourse phenomena) based on a given machine translation task.
 - ▶ For instance, WMT2019 not only provided ratings for each system output but also some manual evaluation and analysis of the outputs, like the English → Czech test suite on coherence by [Rysov et al., 2019].
- ▶ **Specialized test sets** are like normal MT test sets but consist of consecutive sentences that are more densely populated with specific discourse phenomena. Translations are evaluated on such tests sets by means of average quality metrics like BLEU.
 - ▶ E.g. [Voita et al., 2018] build a specialized English → Russian test set by retrieving from OpenSubtitles2016 all the sentences containing pronouns that are coreferent to an expression in the previous sentence.
- ▶ **Contrastive test suites** allow to evaluate MT systems on their ability to rank correct translations higher than the incorrect ones.

Contrastive Test Suites

- **Deixis, Ellipsis, and Lexical Cohesion** [Voita et al., 2019b]

Contrastive Test Suites

- **Deixis, Ellipsis, and Lexical Cohesion** [Voita et al., 2019b]
- **Large Contrastive Test-suite for Pronoun Translation** [Mller et al., 2018]

Contrastive Test Suites

- Deixis, Ellipsis, and Lexical Cohesion [Voita et al., 2019b]
- Large Contrastive Test-suite for Pronoun Translation [Miller et al., 2018]
- **Pronominal Anaphora, Lexical Coherence and Cohesion** [Bawden et al., 2018]

Pronominal Anaphora, Lexical Coherence and Cohesion [[Bawden et al., 2018](#)]

Pronominal Anaphora, Lexical Coherence and Cohesion [[Bawden et al., 2018](#)]

- *Language* English → French (OpenSubtitles2016).

Pronominal Anaphora, Lexical Coherence and Cohesion [Bawden et al., 2018]

- *Language* English → French (OpenSubtitles2016).
- One test suite on **pronominal anaphora** comprised of 50 blocks.

Contrastive Test Suites

Source:

context: Oh, I hate **flies**. Look, there's another one!

current sent.: Don't worry, I'll kill **it** for you.

Target:

- | | | |
|---|---------------|--|
| 1 | context: | Ô je déteste les mouches . Regarde, il y en a une autre ! |
| | correct: | T'inquiète, je la tuerai pour toi. |
| | incorrect: | T'inquiète, je le tuerai pour toi. |
| | | |
| 2 | context: | Ô je déteste les mouchérons . Regarde, il y en a un autre ! |
| | correct: | T'inquiète, je le tuerai pour toi. |
| | incorrect: | T'inquiète, je la tuerai pour toi. |
| | | |
| 3 | context: | Ô je déteste les araignées . Regarde, il y en a une autre ! |
| | semi-correct: | T'inquiète, je la tuerai pour toi. |
| | incorrect: | T'inquiète, je le tuerai pour toi. |
| | | |
| 4 | context: | Ô je déteste les papillons . Regarde, il y en a un autre ! |
| | semi-correct: | T'inquiète, je le tuerai pour toi. |
| | incorrect: | T'inquiète, je la tuerai pour toi. |

Figure: Example block of the pronomial anaphora test suite.

Pronominal Anaphora, Lexical Coherence and Cohesion [Bawden et al., 2018]

- *Language* English → French (OpenSubtitles2016).
- One test suite on **pronominal anaphora** comprised of 50 blocks.
- One on **lexical coherence and cohesion**, comprised of 100 blocks.

Contrastive Test Suites

Source:

context: So what do you say to £50?

current sent.: It's a little **steeper** than I was expecting.

Target:

context: Qu'est-ce que vous en pensez de 50£ ?

correct: C'est un peu plus **cher** que ce que je pensais.

incorrect: C'est un peu plus **raide** que ce que je pensais.

Source:

context: How are your feet holding up?

current sent.: It's a little **steeper** than I was expecting.

Target:

context: Comment vont tes pieds ?

correct: C'est un peu plus **raide** que ce que je pensais.

incorrect: C'est un peu plus **cher** que ce que je pensais.

Figure: Example block of the lexical coherence and cohesion test suite.

Plan

1. Evaluation

1.2 Automatic metrics

1.3 Test Suites

1.4 Remarks

2. Approaches to DNMT

2.5 Concatenation Approaches

2.6 Separate Encoding Approaches

2.7 Cache Approaches

2.8 Exploiting Document-level Monolingual Corpora

2.9 Others

2.10 Remarks and conclusions

Automatic Metrics

Automatic Metrics

- + They are **unexpensive** w.r.t. human annotation.

Automatic Metrics

- + They are **unexpensive** w.r.t. human annotation.
- + they can be easily extended to all languages.

Automatic Metrics

- + They are **unexpensive** w.r.t. human annotation.
- + they can be easily extended to all languages.
- They are **noisy** because they often rely on other imperfect NLP systems. E.g. alignment and coreference systems.

Automatic Metrics

- + They are **unexpensive** w.r.t. human annotation.
- + they can be easily extended to all languages.
- They are **noisy** because they often rely on other imperfect NLP systems. E.g. alignment and coreference systems.
- They might **not be enough correlated with human judgment**:

Automatic Metrics

- + They are **unexpensive** w.r.t. human annotation.
- + they can be easily extended to all languages.
- They are **noisy** because they often rely on other imperfect NLP systems. E.g. alignment and coreference systems.
- They might **not be enough correlated with human judgment**:
 - is the case of APT, for example, which has been shown by [Guillou and Hardmeier, 2018] not to be suitable to evaluate the translation of pronouns with certain functions.

Automatic Metrics

- + They are **unexpensive** w.r.t. human annotation.
- + they can be easily extended to all languages.
- They are **noisy** because they often rely on other imperfect NLP systems. E.g. alignment and coreference systems.
- They might **not be enough correlated with human judgment**:
 - is the case of APT, for example, which has been shown by [Guillou and Hardmeier, 2018] not to be suitable to evaluate the translation of pronouns with certain functions.
- No existing metrics for coherence although it's very relevant for users.

Test Suites

²During scoring, the model is also provided with reference translations as target context (easier). Instead, during translation, the model needs to predict the full sequence, thus being subject to beam search failures and error propagation.

Test Suites

- + They can evaluate discourse phenomena translations with **high precision** and, if well designed, **high recall**.

²During scoring, the model is also provided with reference translations as target context (easier). Instead, during translation, the model needs to predict the full sequence, thus being subject to beam search failures and error propagation.

Test Suites

- + They can evaluate discourse phenomena translations with **high precision** and, if well designed, **high recall**.
- Excepts for specialized test sets (slide 17), test suites have a **limited scope**: fixed language pair, fixed number of context sentences (past and future).

²During scoring, the model is also provided with reference translations as target context (easier). Instead, during translation, the model needs to predict the full sequence, thus being subject to beam search failures and error propagation.

Test Suites

- + They can evaluate discourse phenomena translations with **high precision** and, if well designed, **high recall**.
- Excepts for specialized test sets (slide 17), test suites have a **limited scope**: fixed language pair, fixed number of context sentences (past and future).
- **Contrastive evaluation has limited guarantees**: only permits to conclude whether or not the reference translation is more probable than a contrastive variant. It is not guaranteed at all that the MT system will output such reference translation.²

²During scoring, the model is also provided with reference translations as target context (easier). Instead, during translation, the model needs to predict the full sequence, thus being subject to beam search failures and error propagation.

Plan

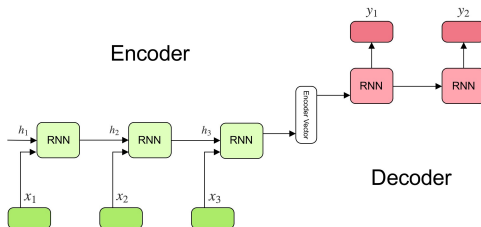
- 1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
- 2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

DNMT architectures are based on traditional encoder-decoder models:

DNMT architectures

DNMT architectures are based on traditional encoder-decoder models:

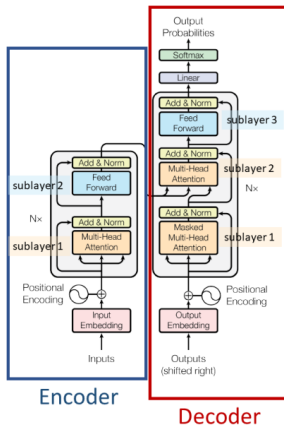
RNN-based (until 2017)



DNMT architectures

DNMT architectures are based on traditional encoder-decoder models:

Transformer-based (afterwards)

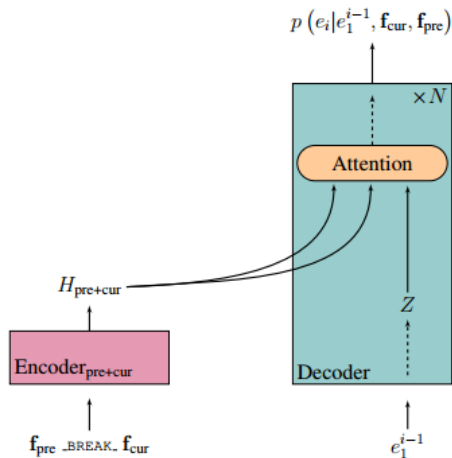


Plan

1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

Concatenation Approaches

Concatenation approaches to DNMT consist in feeding a standard encoder-decoder architecture with a concatenation of sentences.



Concatenation Approaches

For instance:

Concatenation Approaches

For instance:

- ▶ [Tiedemann and Scherrer, 2017] firstly introduced this approach proposing an **RNN-based** model that incorporate the preceding sentence by prepending it to the current one, separated by a <CONCAT> token. They propose two methods:

Concatenation Approaches

For instance:

- [Tiedemann and Scherrer, 2017] firstly introduced this approach proposing an **RNN-based** model that incorporate the preceding sentence by prepending it to the current one, separated by a <CONCAT> token. They propose two methods:
 - **2-TO-2**: the previous and the current sentences are translated together. The translation of the current sentence is then obtained by only retaining the tokens following the concatenation token.

Concatenation Approaches

For instance:

- [Tiedemann and Scherrer, 2017] firstly introduced this approach proposing an **RNN-based** model that incorporate the preceding sentence by prepending it to the current one, separated by a <CONCAT> token. They propose two methods:
 - **2-TO-2**: the previous and the current sentences are translated together. The translation of the current sentence is then obtained by only retaining the tokens following the concatenation token.
 - **2-TO-1**: only the current sentence is translated.

Concatenation Approaches

For instance:

- ▶ [Tiedemann and Scherrer, 2017] firstly introduced this approach proposing an **RNN-based** model that incorporate the preceding sentence by prepending it to the current one, separated by a <CONCAT> token. They propose two methods:
 - ▶ **2-TO-2**: the previous and the current sentences are translated together. The translation of the current sentence is then obtained by only retaining the tokens following the concatenation token.
 - ▶ **2-TO-1**: only the current sentence is translated.
- ▶ [Agrawal et al., 2018, Scherrer et al., 2019] investigated the concatenation approach with the **Transformer** as base model, extending the number of context sentences both on the **source (s:-3,+1)** and the **target (t:-2)** side.

Plan

- 1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
- 2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

Separate Encoding Approaches

Separate encoding approaches to DNMT consist in encoder-decoder models that encode the current and context sentences separately. This can be undertaken by:

Separate Encoding Approaches

Separate encoding approaches to DNMT consist in encoder-decoder models that encode the current and context sentences separately. This can be undertaken by:

- ▶ **Multiple encoders** working in parallel for the current and previous sentence. E.g. [Wang et al., 2017].

Separate Encoding Approaches

Separate encoding approaches to DNMT consist in encoder-decoder models that encode the current and context sentences separately. This can be undertaken by:

- ▶ **Multiple encoders** working in parallel for the current and previous sentence. E.g. [Wang et al., 2017].
- ▶ **Multiple encoders with shared weights.** In this case, the parallel-working encoders not only have the same architecture, but also the same weights. E.g. [Voita et al., 2018].

Separate Encoding Approaches

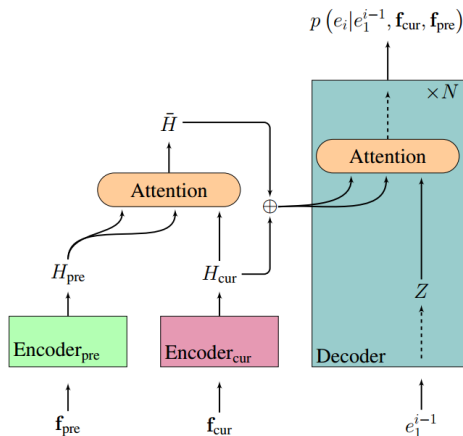
Separate encoding approaches to DNMT consist in encoder-decoder models that encode the current and context sentences separately. This can be undertaken by:

- ▶ **Multiple encoders** working in parallel for the current and previous sentence. E.g. [Wang et al., 2017].
- ▶ **Multiple encoders with shared weights**. In this case, the parallel-working encoders not only have the same architecture, but also the same weights. E.g. [Voita et al., 2018].
- ▶ **Two-pass approaches**, in which the encoder makes a first sentence-level encoding pass of the source, and a second in which it encodes contextual information too. See Slide 38.

Separate Encoding Approaches

Once the encoding of the current and the context sentences has been carried out, they can be integrated in different ways:

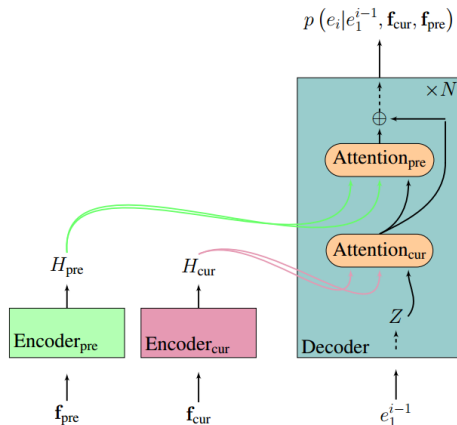
- **Outside** the decoder.
 - (+) symbol represents a gate, a sum or a concatenation.



Separate Encoding Approaches

Once the encoding of the current and the context sentences has been carried out, they can be integrated in different ways:

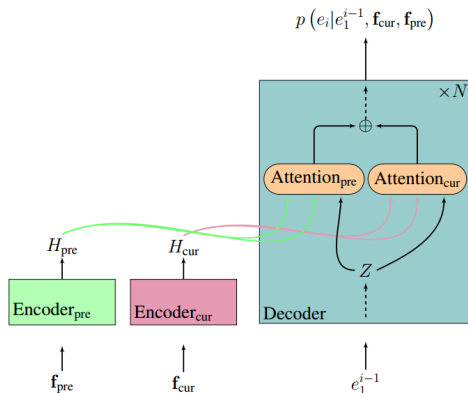
- ▶ **Outside** the decoder.
 - ▶ (+) symbol represents a gate, a sum or a concatenation.
- ▶ **Inside** the decoder, **sequentially**.



Separate Encoding Approaches

Once the encoding of the current and the context sentences has been carried out, they can be integrated in different ways:

- ▶ **Outside** the decoder.
 - ▶ (+) symbol represents a gate, a sum or a concatenation.
- ▶ **Inside** the decoder, **sequentially**.
- ▶ **Inside** the decoder, **in parallel**.



Separate Encoding Approaches

Including target-side context

Despite some have considered including past target-side context harmful because of the *error propagation* problem [Zhang et al., 2018], most recent works have showed it to be of utmost importance for making the most out of context. Past works have successfully included target-side context information in different ways:

Separate Encoding Approaches

Including target-side context

Despite some have considered including past target-side context harmful because of the *error propagation* problem [Zhang et al., 2018], most recent works have showed it to be of utmost importance for making the most out of context. Past works have successfully included target-side context information in different ways:

- Translating past sentences (usually 1) along with the current one, and then discarding them, as in concatenation approaches [Bawden et al., 2018].

Separate Encoding Approaches

Including target-side context

Despite some have considered including past target-side context harmful because of the *error propagation* problem [Zhang et al., 2018], most recent works have showed it to be of utmost importance for making the most out of context. Past works have successfully included target-side context information in different ways:

- ▶ Translating past sentences (usually 1) along with the current one, and then discarding them, as in concatenation approaches [Bawden et al., 2018].
- ▶ By making the decoder attend the target-side hidden representations or embeddings of previously decoded sentences [Miculicich et al., 2018, Voita et al., 2019b, Maruf et al., 2019, Zheng et al., 2020].

Separate Encoding Approaches

Reference	Context	Two-Pass Approach	Outside Integr.	Inside Integr.	Lang. Pair
[Wang et al., 2017]	s:-3		aut...	...aut	Zh→En
[Voita et al., 2018]	s:-1		yes		En→Ru
[Zhang et al., 2018]	s:-2		yes	sequential	Zh→En
[Miculicich et al., 2018]	s:-3; t:-3		yes		Zh/Es→En
[Maruf et al., 2019]	s:all; t:all	optional	yes		En→De
[Zheng et al., 2020]	s:all; t:all	yes	yes		Zh/En→En/De
[Jean et al., 2017]	s:-1			parallel	En→De/Fr
[Bawden et al., 2018]	s:-1; t:-1			parallel	En→Fr
[Fu et al., 2019]	s:all	yes		parallel	En/Zh→De/En
[Voita et al., 2019b]	s:-3; t:-3	yes		parallel*	En→Ru
[Tan et al., 2019]	s:all	yes		parallel	Zh/De→En
[Wang et al., 2019]	s:2			sequential	Fr→En

Positional Embedding Schema

For many approaches to DNMT, the standard positional encoding proposed by [Vaswani et al., 2017] is insufficient because the DNMT system needs to tell context sentences from the current one. For this reason, many strategies have been proposed in the literature, such as:

Positional Embedding Schema

For many approaches to DNMT, the standard positional encoding proposed by [Vaswani et al., 2017] is insufficient because the DNMT system needs to tell context sentences from the current one. For this reason, many strategies have been proposed in the literature, such as:

1. Adding a **sentence distance embedding** to context sentences, that tell the model how far away they are from the current sentence [Voita et al., 2019b].

Positional Embedding Schema

For many approaches to DNMT, the standard positional encoding proposed by [Vaswani et al., 2017] is insufficient because the DNMT system needs to tell context sentences from the current one. For this reason, many strategies have been proposed in the literature, such as:

1. Adding a **sentence distance embedding** to context sentences, that tell the model how far away they are from the current sentence [Voita et al., 2019b].
2. Assign **positional embeddings progressively** to the current sentence, then to the previous one, and so on, so that far away sentences have high values of positional embedding [Li et al., 2019].

Positional Embedding Schema

For many approaches to DNMT, the standard positional encoding proposed by [Vaswani et al., 2017] is insufficient because the DNMT system needs to tell context sentences from the current one. For this reason, many strategies have been proposed in the literature, such as:

1. Adding a **sentence distance embedding** to context sentences, that tell the model how far away they are from the current sentence [Voita et al., 2019b].
2. Assign **positional embeddings progressively** to the current sentence, then to the previous one, and so on, so that far away sentences have high values of positional embedding [Li et al., 2019].
3. Adding a **segment embedding**, similar to classical positional encoding but for the position of the sentence/segment within the document [Zheng et al., 2020].

Plan

1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

Cache Approaches

Cache approaches to DNMT consist in encoder-decoder models that are equipped with one or more caches that store context information. The information stored can belong to both **source side or target side, past and future**.

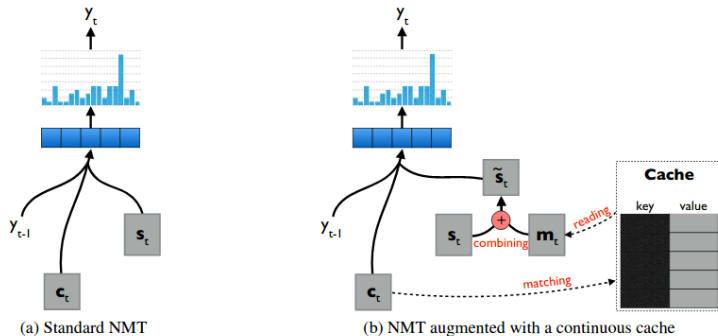
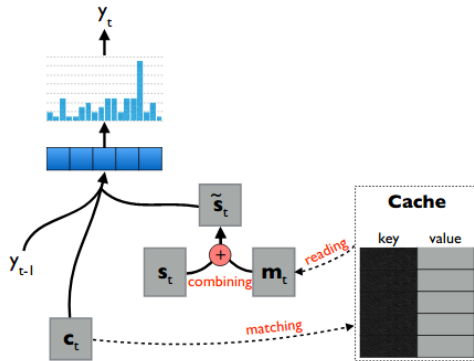


Figure: Continuous cache by [Tu et al., 2017]

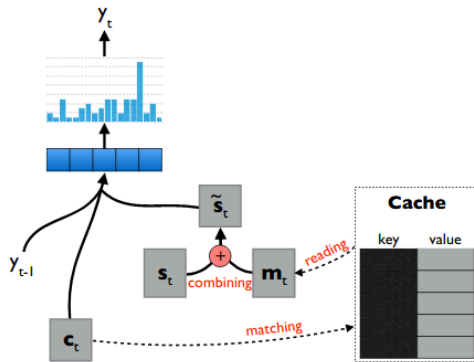
Cache Approaches

Every cache slot is a **key-value** tuple.
With these variables, we can **read** or **write** caches.



Cache Approaches

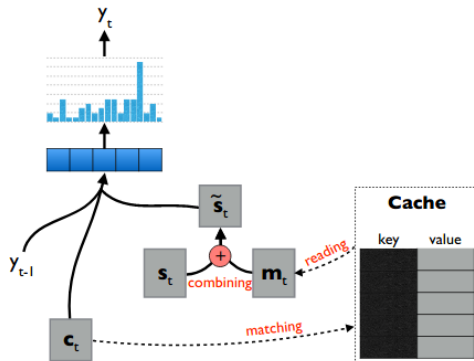
Cache reading involves:



Cache Approaches

Cache reading involves:

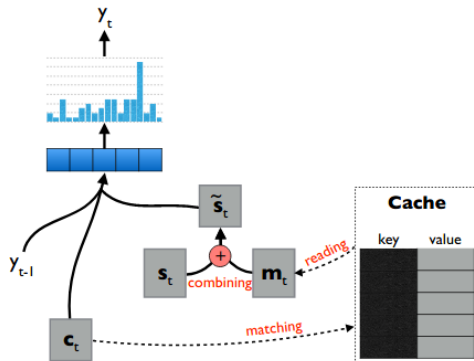
- Soft key matching



Cache Approaches

Cache reading involves:

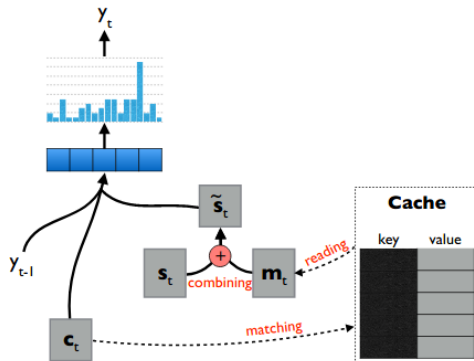
- Soft key matching
- Value reading



Cache Approaches

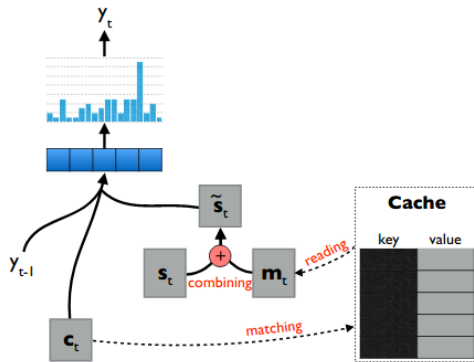
Cache reading involves:

- Soft key matching
- Value reading
- **Combining**



Cache Approaches

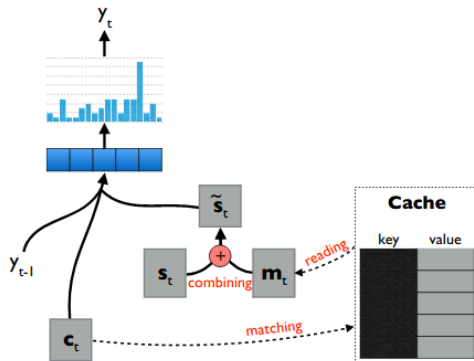
Cache writing can be undertaken after having translated one or more sentences.
For every triplet:



Cache Approaches

Cache writing can be undertaken after having translated one or more sentences. For every triplet:

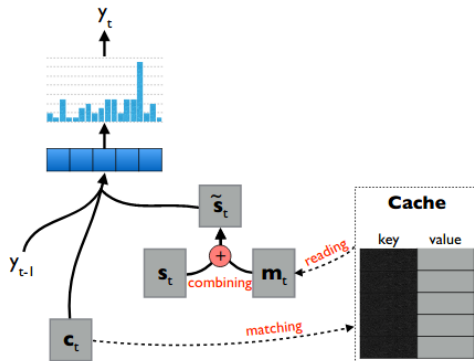
- ▶ If the **key** already exists in the cache, we just update its value.



Cache Approaches

Cache writing can be undertaken after having translated one or more sentences. For every triplet:

- ▶ If the **key** already exists in the cache, we just update its value.
- ▶ Else, we write the key-value tuple in an empty slot, after having emptied the oldest slot if the cache is full.



Cache Approaches

Reference	Caches	Size	Key (Indic.)	Value	Lang. Pair
[Tu et al., 2017]	single	≤ 500	$c_t(y_{k < t})$	$s_{k < t}$	Zh \rightarrow En
[Kuang et al., 2018]	dynamic topic	100 200	c_t	$y_{k < t}$ topic emb.	Zh \rightarrow En
[Maruf and Haffari, 2018]	source target	doc.size	h_t s_t	<i>sent.emb.</i> $s_{k < t}$	Fr/De/Et \rightarrow En

Plan

- 1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
- 2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

On Parallel Corpora for Training

DNMT systems require training on **document-level parallel corpora**. These corpora are usually released during workshops on machine translation like IWSLT and WMT, and hosted on open source web inventories. The most common ones, are extracted from:

On Parallel Corpora for Training

DNMT systems require training on **document-level parallel corpora**. These corpora are usually released during workshops on machine translation like IWSLT and WMT, and hosted on open source web inventories. The most common ones, are extracted from:

- **Movie subtitles (OpenSubtitles)**

On Parallel Corpora for Training

DNMT systems require training on **document-level parallel corpora**. These corpora are usually released during workshops on machine translation like IWSLT and WMT, and hosted on open source web inventories. The most common ones, are extracted from:

- Movie subtitles (OpenSubtitles)
- TED talks (WIT3)

On Parallel Corpora for Training

DNMT systems require training on **document-level parallel corpora**. These corpora are usually released during workshops on machine translation like IWSLT and WMT, and hosted on open source web inventories. The most common ones, are extracted from:

- Movie subtitles (OpenSubtitles)
- TED talks (WIT3)
- News articles (LDC)

On Parallel Corpora for Training

DNMT systems require training on **document-level parallel corpora**. These corpora are usually released during workshops on machine translation like IWSLT and WMT, and hosted on open source web inventories. The most common ones, are extracted from:

- Movie subtitles (OpenSubtitles)
- TED talks (WIT3)
- News articles (LDC)
- Parliamentary interventions (Europarl)

On Parallel Corpora for Training

On Parallel Corpora for Training

- ▶ Unfortunately, document-level parallel **corpora are often insufficient** to train DNMT systems from scratch, although it is often possible to make them converge to a local optimum.

On Parallel Corpora for Training

- ▶ Unfortunately, document-level parallel **corpora are often insufficient** to train DNMT systems from scratch, although it is often possible to make them converge to a local optimum.
- ▶ [Kim et al., 2019, Li et al., 2020] pointed out that when constraining training on such small datasets, model comparison becomes misleading because gains in performance are mainly related to better regularization.

On Parallel Corpora for Training

- ▶ Unfortunately, document-level parallel **corpora are often insufficient** to train DNMT systems from scratch, although it is often possible to make them converge to a local optimum.
- ▶ [Kim et al., 2019, Li et al., 2020] pointed out that when constraining training on such small datasets, model comparison becomes misleading because gains in performance are mainly related to better regularization.
- ▶ A popular solution to this problem is the **two-step training strategy**: [Tu et al., 2017, Zhang et al., 2018, Miculicich et al., 2018]

On Parallel Corpora for Training

- ▶ Unfortunately, document-level parallel **corpora are often insufficient** to train DNMT systems from scratch, although it is often possible to make them converge to a local optimum.
- ▶ [Kim et al., 2019, Li et al., 2020] pointed out that when constraining training on such small datasets, model comparison becomes misleading because gains in performance are mainly related to better regularization.
- ▶ A popular solution to this problem is the **two-step training strategy**: [Tu et al., 2017, Zhang et al., 2018, Miculicich et al., 2018]
 1. Distinguish two integrated components in your model with params $\Theta = [\theta_S; \theta_D]$:

On Parallel Corpora for Training

- ▶ Unfortunately, document-level parallel **corpora are often insufficient** to train DNMT systems from scratch, although it is often possible to make them converge to a local optimum.
- ▶ [Kim et al., 2019, Li et al., 2020] pointed out that when constraining training on such small datasets, model comparison becomes misleading because gains in performance are mainly related to better regularization.
- ▶ A popular solution to this problem is the **two-step training strategy**: [Tu et al., 2017, Zhang et al., 2018, Miculicich et al., 2018]
 1. Distinguish two integrated components in your model with params $\Theta = [\theta_S; \theta_D]$:
 - ▶ A self-standing sentence-level NMT system with parameters θ_S .

On Parallel Corpora for Training

- ▶ Unfortunately, document-level parallel **corpora are often insufficient** to train DNMT systems from scratch, although it is often possible to make them converge to a local optimum.
- ▶ [Kim et al., 2019, Li et al., 2020] pointed out that when constraining training on such small datasets, model comparison becomes misleading because gains in performance are mainly related to better regularization.
- ▶ A popular solution to this problem is the **two-step training strategy**: [Tu et al., 2017, Zhang et al., 2018, Miculicich et al., 2018]
 1. Distinguish two integrated components in your model with params $\Theta = [\theta_S; \theta_D]$:
 - ▶ A self-standing sentence-level NMT system with parameters θ_S .
 - ▶ Some context-handling modules with parameters θ_D .

On Parallel Corpora for Training

- ▶ Unfortunately, document-level parallel **corpora are often insufficient** to train DNMT systems from scratch, although it is often possible to make them converge to a local optimum.
- ▶ [Kim et al., 2019, Li et al., 2020] pointed out that when constraining training on such small datasets, model comparison becomes misleading because gains in performance are mainly related to better regularization.
- ▶ A popular solution to this problem is the **two-step training strategy**: [Tu et al., 2017, Zhang et al., 2018, Miculicich et al., 2018]
 1. Distinguish two integrated components in your model with params $\Theta = [\theta_S; \theta_D]$:
 - ▶ A self-standing sentence-level NMT system with parameters θ_S .
 - ▶ Some context-handling modules with parameters θ_D .
 2. Train θ_S independently on a sentence-level parallel corpus C_S .

On Parallel Corpora for Training

- ▶ Unfortunately, document-level parallel **corpora are often insufficient** to train DNMT systems from scratch, although it is often possible to make them converge to a local optimum.
- ▶ [Kim et al., 2019, Li et al., 2020] pointed out that when constraining training on such small datasets, model comparison becomes misleading because gains in performance are mainly related to better regularization.
- ▶ A popular solution to this problem is the **two-step training strategy**: [Tu et al., 2017, Zhang et al., 2018, Miculicich et al., 2018]
 1. Distinguish two integrated components in your model with params $\Theta = [\theta_S; \theta_D]$:
 - ▶ A self-standing sentence-level NMT system with parameters θ_S .
 - ▶ Some context-handling modules with parameters θ_D .
 2. Train θ_S independently on a sentence-level parallel corpus C_S .
 3. Train θ_D on a document-level parallel corpus C_D while fine-tuning θ_S , or freezing them [Zhang et al., 2018].

Exploiting Document-level Monolingual Corpora

Another solution to the lack of vast document-level parallel corpora is leveraging on huge *monolingual* document-level corpora like BookCorpus [Zhu et al., 2015] and PG-19 [Rae et al., 2019]. In the literature, we can find various approaches to leverage monolingual corpora:

Exploiting Document-level Monolingual Corpora

Another solution to the lack of vast document-level parallel corpora is leveraging on huge *monolingual* document-level corpora like BookCorpus [Zhu et al., 2015] and PG-19 [Rae et al., 2019]. In the literature, we can find various approaches to leverage monolingual corpora:

- **Back-translate** target-side corpus to augment dl corpus [Sugiyama and Yoshinaga, 2019].

Exploiting Document-level Monolingual Corpora

Another solution to the lack of vast document-level parallel corpora is leveraging on huge *monolingual* document-level corpora like BookCorpus [Zhu et al., 2015] and PG-19 [Rae et al., 2019]. In the literature, we can find various approaches to leverage monolingual corpora:

- **Back-translate** target-side corpus to augment dl corpus [Sugiyama and Yoshinaga, 2019].
- Train **context-aware language models** on target/source-side corpus, then:

Exploiting Document-level Monolingual Corpora

Another solution to the lack of vast document-level parallel corpora is leveraging on huge *monolingual* document-level corpora like BookCorpus [Zhu et al., 2015] and PG-19 [Rae et al., 2019]. In the literature, we can find various approaches to leverage monolingual corpora:

- ▶ **Back-translate** target-side corpus to augment dl corpus [Sugiyama and Yoshinaga, 2019].
- ▶ Train **context-aware language models** on target/source-side corpus, then:
 - ▶ Generate translations by fusing the decoder and the LM's scores to candidate words [Martnez Garcia et al., 2019].

Exploiting Document-level Monolingual Corpora

Another solution to the lack of vast document-level parallel corpora is leveraging on huge *monolingual* document-level corpora like BookCorpus [Zhu et al., 2015] and PG-19 [Rae et al., 2019]. In the literature, we can find various approaches to leverage monolingual corpora:

- ▶ **Back-translate** target-side corpus to augment dl corpus [Sugiyama and Yoshinaga, 2019].
- ▶ Train **context-aware language models** on target/source-side corpus, then:
 - ▶ Generate translations by fusing the decoder and the LM's scores to candidate words [Martnez Garcia et al., 2019].
 - ▶ Initialize the encoder (or decoder) of a DNMT model [Li et al., 2019].

Exploiting Document-level Monolingual Corpora

Another solution to the lack of vast document-level parallel corpora is leveraging on huge *monolingual* document-level corpora like BookCorpus [Zhu et al., 2015] and PG-19 [Rae et al., 2019]. In the literature, we can find various approaches to leverage monolingual corpora:

- ▶ **Back-translate** target-side corpus to augment dl corpus [Sugiyama and Yoshinaga, 2019].
- ▶ Train **context-aware language models** on target/source-side corpus, then:
 - ▶ Generate translations by fusing the decoder and the LM's scores to candidate words [Martnez Garcia et al., 2019].
 - ▶ Initialize the econdor (or decoder) of a DNMT model [Li et al., 2019].
- ▶ Train **Automatic Post Editing** systems on target-side corpus (See next slide).

Automatic Post Editing (APE)

[[Voita et al., 2019a](#)] devised an APE system called DocRepair, that turns a sentence-level translation into a context-aware translation. DocRepair can work on top of whatever sentence-level MT system.

DocRepair

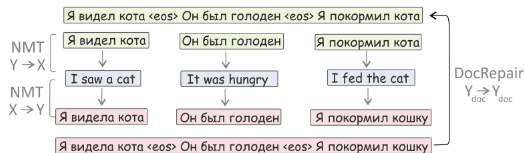


Figure 1: **Training procedure of DocRepair**. First, round-trip translations of individual sentences are produced to form an inconsistent text fragment (in the example, both genders of the speaker and the cat became inconsistent). Then, a repair model is trained to produce an original text from the inconsistent one.



Figure 2: The process of producing document-level translations at **test time** is two-step: (1) sentences are translated independently using a sentence-level model, (2) DocRepair model corrects translation of the resulting text fragment.

Plan

1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

Approaches Including Additional Discourse Information as Input

These approaches consist in concatenation approaches or separate encoding approaches that also integrate discourse-related information as additional input features. Examples of extra features are:

Approaches Including Additional Discourse Information as Input

These approaches consist in concatenation approaches or separate encoding approaches that also integrate discourse-related information as additional input features. Examples of extra features are:

- ▶ Lexical chains of semantically similar words to promote word sense disambiguation [Rios Gonzales et al., 2017].

Approaches Including Additional Discourse Information as Input

These approaches consist in concatenation approaches or separate encoding approaches that also integrate discourse-related information as additional input features. Examples of extra features are:

- Lexical chains of semantically similar words to promote word sense disambiguation [Rios Gonzales et al., 2017].
- Coreference chains to promote coreference resolution [Stojanovski and Fraser, 2018, Ohtani et al., 2019].

Learning Approaches

[[Jean and Cho, 2019](#)] looked at the problem from a learning perspective and designed a regularisation term to encourage a DNMT model to exploit the additional context in a useful way . This regularisation term is applied at the token, sentence and corpus levels and is based on pair-wise ranking loss, that is, it helps to assign a higher log-probability to a translation paired with the correct context than to the translation without context.

Plan

1. Evaluation
 - 1.2 Automatic metrics
 - 1.3 Test Suites
 - 1.4 Remarks
2. Approaches to DNMT
 - 2.5 Concatenation Approaches
 - 2.6 Separate Encoding Approaches
 - 2.7 Cache Approaches
 - 2.8 Exploiting Document-level Monolingual Corpora
 - 2.9 Others
 - 2.10 Remarks and conclusions

Possible Future Research Directions

Possible Future Research Directions

- Build a large DL corpus for training systems, or find automatic approaches to generate synthetic data other than back-translation.

Possible Future Research Directions

- Build a large DL corpus for training systems, or find automatic approaches to generate synthetic data other than back-translation.
 - E.g. imputing context sentences [Jean et al., 2019].

Possible Future Research Directions

- Build a large DL corpus for training systems, or find automatic approaches to generate synthetic data other than back-translation.
 - E.g. imputing context sentences [[Jean et al., 2019](#)].
- Design models exploiting full context in a memory-efficient way:

Possible Future Research Directions

- Build a large DL corpus for training systems, or find automatic approaches to generate synthetic data other than back-translation.
 - E.g. imputing context sentences [[Jean et al., 2019](#)].
- Design models exploiting full context in a memory-efficient way:
 - **Dynamic context integration.**

Possible Future Research Directions

- Build a large DL corpus for training systems, or find automatic approaches to generate synthetic data other than back-translation.
 - E.g. imputing context sentences [[Jean et al., 2019](#)].
- Design models exploiting full context in a memory-efficient way:
 - Dynamic context integration.
 - **Caches integrated to Transformer-based models.**

Possible Future Research Directions

- ▶ Build a large DL corpus for training systems, or find automatic approaches to generate synthetic data other than back-translation.
 - ▶ E.g. imputing context sentences [[Jean et al., 2019](#)].
- ▶ Design models exploiting full context in a memory-efficient way:
 - ▶ Dynamic context integration.
 - ▶ Caches integrated to Transformer-based models.
- ▶ Design automatic post-processing models that are lightweight and can be trained on little data [[Kim et al., 2019](#)].

Possible Future Research Directions



- ▶ Build a large DL corpus for training systems, or find automatic approaches to generate synthetic data other than back-translation.
 - ▶ E.g. imputing context sentences [[Jean et al., 2019](#)].
- ▶ Design models exploiting full context in a memory-efficient way:
 - ▶ Dynamic context integration.
 - ▶ Caches integrated to Transformer-based models.
- ▶ Design automatic post-processing models that are lightweight and can be trained on little data [[Kim et al., 2019](#)].
- ▶ Study pre-trained language models for DNMT decoder.

Possible Future Research Directions



- Build a large DL corpus for training systems, or find automatic approaches to generate synthetic data other than back-translation.
 - E.g. imputing context sentences [Jean et al., 2019].
- Design models exploiting full context in a memory-efficient way:
 - Dynamic context integration.
 - Caches integrated to Transformer-based models.
- Design automatic post-processing models that are lightweight and can be trained on little data [Kim et al., 2019].
- Study pre-trained language models for DNMT decoder.
- Study other learning methods that foster document-level modeling [Jean and Cho, 2019].

Thank you for your attention!

References I

-  Agrawal, R. R., Turchi, M., and Negri, M. (2018).
Contextual Handling in Neural Machine Translation: Look Behind, Ahead and on
Both Sides.
pages 11–20.
00007 Accepted: 2018-08-08T15:15:28Z.
-  Bawden, R., Sennrich, R., Birch, A., and Haddow, B. (2018).
Evaluating Discourse Phenomena in Neural Machine Translation.
*In Proceedings of the 2018 Conference of the North American Chapter of the
Association for Computational Linguistics: Human Language Technologies,
Volume 1 (Long Papers)*, pages 1304–1313, New Orleans, Louisiana. Association
for Computational Linguistics.
00056.




References II

-  Fellbaum, C. (1998).
A Semantic Network of English: The Mother of All WordNets.
Computers and the Humanities, 32(2):209–220.
00194.
-  Fu, H., Liu, C., and Sun, J. (2019).
Reference Network for Neural Machine Translation.
In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 3002–3012, Florence, Italy. Association for Computational Linguistics.
00000.



References III

-  Guillo, L. and Hardmeier, C. (2018).
Automatic Reference-Based Evaluation of Pronoun Translation Misses the Point.
In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4797–4802, Brussels, Belgium. Association for Computational Linguistics.
00008.
-  Hajlaoui, N. and Popescu-Belis, A. (2013).
Assessing the Accuracy of Discourse Connective Translations: Validation of an Automatic Metric.
In Gelbukh, A., editor, *Computational Linguistics and Intelligent Text Processing*, Lecture Notes in Computer Science, pages 236–247, Berlin, Heidelberg. Springer.
00000.


References IV

-  Jean, S., Bapna, A., and Firat, O. (2019).
Fill in the Blanks: Imputing Missing Sentences for Larger-Context Neural Machine Translation.
arXiv:1910.14075 [cs].
00000 arXiv: 1910.14075.
-  Jean, S. and Cho, K. (2019).
Context-Aware Learning for Neural Machine Translation.
arXiv:1903.04715 [cs].
00003 arXiv: 1903.04715.
-  Jean, S., Lauly, S., Firat, O., and Cho, K. (2017).
Does Neural Machine Translation Benefit from Larger Context?
arXiv:1704.05135 [cs, stat].
00039 arXiv: 1704.05135.



References V

-  Jwalapuram, P., Joty, S., Temnikova, I., and Nakov, P. (2019).
Evaluating Pronominal Anaphora in Machine Translation: An Evaluation Measure
and a Test Suite.
arXiv:1909.00131 [cs].
00002 arXiv: 1909.00131.
-  Kim, Y., Tran, D. T., and Ney, H. (2019).
When and Why is Document-level Context Useful in Neural Machine Translation?
arXiv:1910.00294 [cs].
00001 arXiv: 1910.00294.



References VI

-  Kuang, S., Xiong, D., Luo, W., and Zhou, G. (2018).
Modeling Coherence for Neural Machine Translation with Dynamic and Topic
Caches.
In Proceedings of the 27th International Conference on Computational Linguistics,
pages 596–606, Santa Fe, New Mexico, USA. Association for Computational
Linguistics.
00012.
-  Li, B., Liu, H., Wang, Z., Jiang, Y., Xiao, T., Zhu, J., Liu, T., and Li, C. (2020).
Does Multi-Encoder Help? A Case Study on Context-Aware Neural Machine
Translation.
arXiv:2005.03393 [cs].
00000 arXiv: 2005.03393.



References VII

-  Li, L., Jiang, X., and Liu, Q. (2019).
Pretrained Language Models for Document-Level Neural Machine Translation.
arXiv:1911.03110 [cs].
00001 arXiv: 1911.03110.
-  Lubli, S., Sennrich, R., and Volk, M. (2018).
Has Machine Translation Achieved Human Parity? A Case for Document-level
Evaluation.
arXiv:1808.07048 [cs].
00038 arXiv: 1808.07048.



References VIII

-  Martnez Garcia, E., Creus, C., and Espaa-Bonet, C. (2019).
Context-Aware Neural Machine Translation Decoding.
In Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019), pages 13–23, Hong Kong, China. Association for Computational Linguistics.
00000.
-  Maruf, S. and Haffari, G. (2018).
Document Context Neural Machine Translation with Memory Networks.
In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1275–1284, Melbourne, Australia. Association for Computational Linguistics.
00000.




References IX

-  Maruf, S., Martins, A. F. T., and Haffari, G. (2019).
Selective Attention for Context-aware Neural Machine Translation.
arXiv:1903.08788 [cs].
00012.
-  Miculicich, L., Ram, D., Pappas, N., and Henderson, J. (2018).
Document-Level Neural Machine Translation with Hierarchical Attention
Networks.
arXiv:1809.01576 [cs].
00029 arXiv: 1809.01576.



References X

-  Miculicich Werlen, L. and Popescu-Belis, A. (2017).
Validation of an Automatic Metric for the Accuracy of Pronoun Translation (APT).
In Proceedings of the Third Workshop on Discourse in Machine Translation, pages 17–25, Copenhagen, Denmark. Association for Computational Linguistics. 00000.
-  Mller, M., Rios, A., Voita, E., and Sennrich, R. (2018).
A Large-Scale Test Set for the Evaluation of Context-Aware Pronoun Translation in Neural Machine Translation.
In Proceedings of the Third Conference on Machine Translation: Research Papers, pages 61–72, Brussels, Belgium. Association for Computational Linguistics. 00010.

References XI

-  Ohtani, T., Kamigaito, H., Nagata, M., and Okumura, M. (2019).
Context-aware Neural Machine Translation with Coreference Information.
In Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019), pages 45–50, Hong Kong, China. Association for Computational Linguistics.
00000.
-  Porter, M. (1980).
An algorithm for suffix stripping.
Program, 40(3):211–218.
10830.
-  Rae, J. W., Potapenko, A., Jayakumar, S. M., and Lillicrap, T. P. (2019).
Compressive Transformers for Long-Range Sequence Modelling.
arXiv:1911.05507 [cs, stat].
00000 arXiv: 1911.05507.

References XII

-  Rios Gonzales, A., Mascarell, L., and Sennrich, R. (2017).
Improving Word Sense Disambiguation in Neural Machine Translation with Sense
Embeddings.
In Proceedings of the Second Conference on Machine Translation, pages 11–19,
Copenhagen, Denmark. Association for Computational Linguistics.
00030.
-  Rysov, K., Rysov, M., Musil, T., Polkov, L., and Bojar, O. (2019).
A Test Suite and Manual Evaluation of Document-Level NMT at WMT19.
*In Proceedings of the Fourth Conference on Machine Translation (Volume 2:
Shared Task Papers, Day 1)*, pages 455–463, Florence, Italy. Association for
Computational Linguistics.
00001.

References XIII

-  Scherrer, Y., Tiedemann, J., and Loiciga, S. (2019).
Analysing concatenation approaches to document-level NMT in two different domains.
In Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019), pages 51–61, Hong Kong, China. Association for Computational Linguistics.
00002.
-  Stojanovski, D. and Fraser, A. (2018).
Coreference and Coherence in Neural Machine Translation: A Study Using Oracle Experiments.
In Proceedings of the Third Conference on Machine Translation: Research Papers, pages 49–60, Brussels, Belgium. Association for Computational Linguistics.
00003.

References XIV



-  Sugiyama, A. and Yoshinaga, N. (2019).
Data augmentation using back-translation for context-aware neural machine translation.
In Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019), pages 35–44, Hong Kong, China. Association for Computational Linguistics.
00001.
-  Tan, X., Zhang, L., Xiong, D., and Zhou, G. (2019).
Hierarchical Modeling of Global Context for Document-Level Neural Machine Translation.
In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 1576–1585, Hong Kong, China. Association for Computational Linguistics.

00004.

-  Tiedemann, J. and Scherrer, Y. (2017).
Neural Machine Translation with Extended Context.
In *Proceedings of the Third Workshop on Discourse in Machine Translation*, pages
82–92, Copenhagen, Denmark. Association for Computational Linguistics.
00040.

-  Tu, Z., Liu, Y., Shi, S., and Zhang, T. (2017).
Learning to Remember Translation History with a Continuous Cache.
arXiv:1711.09367 [cs].
00041 arXiv: 1711.09367.



References XVI

-  Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017).
Attention Is All You Need.
arXiv:1706.03762 [cs].
05728 arXiv: 1706.03762.
-  Voita, E., Sennrich, R., and Titov, I. (2019a).
Context-Aware Monolingual Repair for Neural Machine Translation.
arXiv:1909.01383 [cs].
00003 arXiv: 1909.01383.



References XVII

-  Voita, E., Sennrich, R., and Titov, I. (2019b).
When a Good Translation is Wrong in Context: Context-Aware Machine Translation Improves on Deixis, Ellipsis, and Lexical Cohesion.
In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 1198–1212, Florence, Italy. Association for Computational Linguistics.
00007.
-  Voita, E., Serdyukov, P., Sennrich, R., and Titov, I. (2018).
Context-Aware Neural Machine Translation Learns Anaphora Resolution.
In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1264–1274, Melbourne, Australia. Association for Computational Linguistics.
00049.

References XVIII

-  Wang, L., Tu, Z., Way, A., and Liu, Q. (2017).
Exploiting Cross-Sentence Context for Neural Machine Translation.
In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2826–2831, Copenhagen, Denmark. Association for Computational Linguistics.
00048.
-  Wang, X., Weston, J., Auli, M., and Jernite, Y. (2019).
Improving Conditioning in Context-Aware Sequence to Sequence Models.
00002.

References XIX


-  Wong, B. T. M. and Kit, C. (2012).
Extending Machine Translation Evaluation Metrics with Lexical Cohesion to Document Level.
In Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, pages 1060–1068, Jeju Island, Korea. Association for Computational Linguistics. 00044.
-  Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., Klingner, J., Shah, A., Johnson, M., Liu, X., Kaiser, ., Gouws, S., Kato, Y., Kudo, T., Kazawa, H., Stevens, K., Kurian, G., Patil, N., Wang, W., Young, C., Smith, J., Riesa, J., Rudnick, A., Vinyals, O., Corrado, G., Hughes, M., and Dean, J. (2016).
Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation.


References XX

arXiv:1609.08144 [cs].

00000 arXiv: 1609.08144.

-  Zhang, J., Luan, H., Sun, M., Zhai, F., Xu, J., Zhang, M., and Liu, Y. (2018). Improving the Transformer Translation Model with Document-Level Context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 533–542, Brussels, Belgium. Association for Computational Linguistics.
00028.

-  Zheng, Z., Yue, X., Huang, S., Chen, J., and Birch, A. (2020). Toward Making the Most of Context in Neural Machine Translation.
arXiv:2002.07982 [cs].
00000 arXiv: 2002.07982.

-  Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., and Fidler, S. (2015).
Aligning Books and Movies: Towards Story-like Visual Explanations by Watching Movies and Reading Books.
arXiv:1506.06724 [cs].
00450 arXiv: 1506.06724.