# Killing Two Flies with One Stone: An Attempt to Break LLMs Using English→Icelandic Idioms and Proper Names

**Bjarki Ármannsson, Hinrik Hafsteinsson, Atli Jasonarson, Steinþór Steingrímsson**
The Árni Magnússon Institute for Icelandic Studies
Reykjavík, Iceland
bjarki.armannsson,hinrik.hafsteinsson,atli.jasonarson,
steinthor.steingrimsson@arnastofnun.is

## Abstract

This paper presents the submission of the Árni Magnússon Institute's team to the WMT24 test suite subtask, focusing on idiomatic expressions and proper names for the English→Icelandic translation direction.

Intuitively and empirically, idioms and proper names are known to be a significant challenge for modern translation models. We create two different test suites. The first evaluates the competency of MT systems in translating common English idiomatic expressions, as well as testing whether systems can distinguish between those expressions and the same phrases when used in a literal context. The second test suite consists of place names that should be translated into their Icelandic exonyms (and correctly inflected) and pairs of Icelandic names that share a surface form between the male and female variants, so that incorrect translations impact meaning as well as readability.

The scores reported are relatively low, especially for idiomatic expressions and place names, and indicate considerable room for improvement.

## 1 Introduction

Significant advances in machine translation have in recent years been achieved by integrating Large Language Models (LLMs) into neural translation systems (Xu et al., 2024). Careful analysis, however, has repeatedly shown that despite recording higher scores and producing text with greater fluency compared to previous state-of-the-art neural systems, the translations produced by LLMs are still far from perfect and can include significant biases, misinformation and hallucinations (Hendy et al., 2023), half-hidden in the impressive-looking output. Aiming to expose "weaknesses and serious flaws" of these systems that might otherwise get "hidden in the average", the theme of this year's WMT test suite subtask is "Help us break LLMs",

with organizers asking for custom test sets focusing on phenomena that can provide specific challenges for LLM-based systems. This paper describes the efforts of the Árni Magnússon Institute's team to pick holes in otherwise seemingly fluent English→Icelandic translations.

We experiment with two main features we believe should prove particularly challenging for English→Icelandic LLM-based machine translation systems; idiomatic expressions and proper names. More specifically, we focus on:

- **Idiomatic expressions in English and their literal counterparts:** In the first of our two test sets, we investigate idiomatic expressions in English which do not directly translate to Icelandic. Where possible, we also include 'inverse' examples of usage in a literal form (as in "Are you supposed to **chew the fat** from steak?" or "Blow into the balloon and **tie the knot** without letting the air out.") to give an idea of the translation models' ability to correctly switch between literal and non-literal translations of the same phrase, depending on context.

- **Proper names:** In our second test set, we also consider names of both people and places. We carefully curate a list of city and area names in English that should be translated to their common Icelandic names (and correctly inflected). We then include a list of simple sentences containing both Icelandic and English given names. For the Icelandic names, we observe whether they are correctly inflected in the Icelandic text (which impacts not only the text's readability, but also its meaning). Common English names, meanwhile, are included to test that the models don't 'translate' them to Icelandic – i.e. alter them in some unintended way.

We release our test suites and evaluation code

for others to build on and to allow for further comparison between future models in these categories.[1]

## 2 Related Work

Idiomatic expressions (and multi-word expressions (MWEs) in general) have been the focus of much work in the field of machine translation in recent years and the construction of impressive idiom datasets has been carried out for many other languages and language pairs. See e.g. Stap et al. (2024) for English↔German and Russian→English, Tang (2022) for Chinese→English, Fadaee et al. (2018) for English↔German and Haagsma et al. (2020) and Adewumi et al. (2022) for monolingual datasets of English idiomatic expressions.

Macketanz et al. (2022) include idioms among many other interesting linguistic phenomena in their dataset for English↔German and English→Russian and we took some inspiration from their work when deciding on our scoring format. Halldórsson et al. (2022) list Icelandic idioms with English equivalents, this dataset is described and discussed in more detail in Steingrímsson et al. (2024). We are not aware of any dataset for the English→Icelandic translation direction published previous to our work.

In recent years, the emergence of LLMs has led to work investigating how they handle the translations of idioms and MWEs compared with previous models. Raunak et al. (2023), using measures of 'literalness', find that GPT models produce less literal translations between English and German, Chinese, and Russian than previous neural models, a difference most pronounced in the case of idiomatic expressions. Finally, Shwartz (2021) provide an accessible overview of the kinds of problems posed by MWEs for language models in general.

## 3 Methodology

### 3.1 Idiomatic Expressions

We make use of the set of potential idiomatic expressions defined in the PIE Corpus (Adewumi et al., 2022) and, for each expression we use, extract two examples of usage from the NewsCrawl corpus of WMT 2023 (Kocmi et al., 2023)[2] For our purposes, we narrow the PIE set down from 591 expressions to 199. Our aim was to remove those we deem too rare or obscure to be truly relevant (e.g. *horses for courses* or *monkey's uncle*) for model comparison, expressions which directly (or more or less directly) translate between English and Icelandic (e.g. "open the floodgates" has an Icelandic equivalent, "opna flóðgáttirnar") and those for which we find no example usage in the NewsCrawl corpus. We make the number of expressions an even 200 by adding one that was not in the PIE corpus: "kill two birds with one stone".

Each of the 400 example sentences - two examples for each of the 200 selected idioms - is then manually reviewed to make sure that the relevant idiomatic expression is being used in the intended, non-literal sense. To further increase the difficulty of the task (though still keeping it trivial for fluent human speakers of Icelandic and English), we also try and test the models on their ability to translate the words in these expressions literally when appropriate. We include 223 additional example sentences, for as many expressions as we were able, where the expression is used in a literal sense (or in a few cases, very slightly altered to try and exploit the likelihood bias of LLMs).[3] These examples are largely taken from the NewsCrawl corpus but synthetic in some cases.

To evaluate the models' performance, we construct two 'positive' sets of Icelandic word forms or multiword expressions for each idiom. One set contains words that we would expect to find in a literal translation of the phrase, the other words or phrases that could be expected to appear in a suitable, non-literal translation of the idiomatic expression. In many cases, we also construct 'negative' sets of words that instantly lead to a sentence being marked incorrect, such as the Icelandic words for "weather" or "pink" for idiomatic translations of the phrases "under the weather" and "in the pink". An Icelandic translation of an example sentence in English is marked as correct if it contains any of the words in the set of 'positive' words (in any lexical form) **and** it contains none of the words in the set of 'negative' words (see Table 1).

---

[3]Early inspiration for this project was provided by the one idiomatic expression we added from outside the PIE corpus: "kill two birds with one stone". We noticed a prominent online translation service correctly translated this to the equivalent Icelandic phrase, "slá tvær flugur í einu höggi" (lit. *hit two flies in one strike*), whereas a phrase like "He killed two birds yesterday" would be wrongly translated as "Hann drap tvær flugur í gær" (lit. *He killed two flies yesterday*), exposing a weakness particular to neural and LLM-based systems. Indeed, four of the systems tested here made this particular mistake.

| Source sentence | Possible translations | Evaluation |
|---|---|---|
| Why Fleabag is **in the pink**! | Fleabag er *í góðum málum*! | ✓(Positive match) |
| | Fleabag er **í bleiku**! | ✗ (Negative match) |
| The young woman **in the pink** continued to throw punches [...] | Unga konan **í góðum málum** lét hnefana tala áfram [...] | ✗ (No positive match) |
| | Unga konan **í bleiku fötunum** lét hnefana tala áfram [...] | ✓(Positive match) |

Table 1: Fabricated example translations into Icelandic of two English sentences containing the phrase "in the pink", both from our test suite. The first English sentence uses the phrase in an idiomatic sense (meaning *in good health* or *in a state of well-being*) and the second seemingly in a literal sense (the full context, not included in the table, is: "before another wades in"). For the idiomatic sentence, we automatically mark it as correct if a match is found from a list of possible Icelandic translations (here the phrase "í góðum málum") *and* no match is found from a list of negative matches (here the lexeme "bleikur", meaning *pink*). For the literal sentence, meanwhile, some form of "bleikur" is required for a correct marking.

During our manual evaluation, we further whittled down our set as we decided a few sentences we had decided to include were actually not testing what they were meant to test (as some were, for instance, more linguistically acceptable when translated directly into Icelandic than we originally felt during the construction of our test set). We removed a total of 25 sentences this way, bringing the total of 'idiomatic' examples in our set to 397 and the total of 'literal' examples to 201. Note that although these examples were removed after we received their translations from the tested models, they are not included in our scoring.

### 3.2 Proper Names

For our testing of place names, we construct our own list of 52 names of cities and areas that we argue would be highly unusual not to translate into their Icelandic names.[4]

As a reference when collecting our place names, we make use of Wikipedia's list of Icelandic exonyms.[5] We use only a small subset of that list, however. Aiming to err on the side of caution, we try to include only place names where native speakers would be in more or less complete agreement to apply their Icelandic names rather than the ones used in English (e.g. the name "Kaupmannahöfn" for Copenhagen is invariably used, whereas "Lundúnir" for London is very rare and mostly used in a colourful or joking manner.[6] We

also leave out cases where the differences between the names used in English and in Icelandic only have to do with pronunciation or minor differences in spelling. In addition to the Icelandic exonyms we select, we make sure to also include several examples of cities where the local name is the one more generally used and English speakers use a rarer (typically French-derived) name (e.g. "München" rather than the English "Munich").

We then construct example sentences in English where each of our selected place names is used in four different contexts, corresponding to each of the four grammatical cases in Icelandic. (The exceptions are "Paris" and "Berlin", which are only tested in the genitive as they are practically the same as in English in the other three cases.) Our motivation is that due to the richer morphology of Icelandic, an accurate translation model needs to be able to map the same lexical form in English to several different forms in Icelandic, depending on the context (and this particular mapping is perhaps a problem better suited to older models than state-of-the-art LLM-based ones).

We try to avoid the possibility that our sentences will be translated into Icelandic in a way that is generally correct but uses a different syntactic structure or wording than we anticipate, which would lend itself to the use of a different grammatical case than the one we intend to test for. We do this by keeping our example sentences short and simple and choose case-governing words and prepositions carefully to maximize the probability of a particular translation in Icelandic (e.g. the sentence "The flight

---

[4]There exist context-dependent exceptions to this, of course, such as the name of a sport club or particular institution from a certain city. Our example sentences, however, refer clearly to the cities in general.

[5]https://en.wikipedia.org/wiki/Icelandic_exonyms

[6]One anonymous reviewer asked whether we had considered incorporating a native speaker survey in order to validate our choices. While the suggestion is certainly a good one, it is beyond the scope of this particular work.

| Source sentence | Possible translations | Evaluation |
|---|---|---|
| **Helgi** dreams of flying | **Helgi** dreymir um að fljúga | ✗ (Ungrammatical) |
| | **Helga** dreymir um að fljúga | ✓ |
| **Helga** dreams of flying | **Helga** dreymir um að fljúga | ✗ (Refers to Helgi, not Helga) |
| | **Helgu** dreymir um að fljúga | ✓ |

Table 2: Examples of possible translations of the phrase "dreams of flying". In Icelandic, the verb "dreyma" (*to dream*) takes a subject argument in the accusative case, which requires a translation system to alter the form of the given name in the English text. Left unaltered in Icelandic, the male name "Helgi" renders the sentence ungrammatical and the female name "Helga" would cause the reader to interpret the sentence to refer to a male called Helgi instead.

from Tórshavn to Gothenburg was delayed until the morning" should almost certainly be translated using the prepositions "frá" and "til" for "from" and "to", governing the dative and genitive cases respectively.)

Given names, both in Icelandic and English, constitute the final part of our test suite. As in the case of the place names, we construct simple sentences in English containing Icelandic names and meant to test for each of the four grammatical cases. For this task, we chose a specific subset of common names in Icelandic: male-female pairings that take the weak inflection, e.g. "Helgi"-"Helga" and "Gunni"-"Gunna", where the male name has the ending -"i" in the nominative case but -"a" in oblique cases and the female name has the ending -"a" in the nominative but -"u" in the oblique cases (and possibly also a u-umlaut as in "Svala" → "Svölu").

These name pairs, of which we select 45 from the Database of Icelandic Morphology (Bjarnadóttir et al., 2019),[7] are chosen as they seem to present a particular challenge for translation systems compared to names that take the strong declension. In constructing our test suite, we found that available models seemed to perform at random when asked to translate sentences containing these names in different cases, presumably due to the ambiguity of the lexical forms ending in -"a", which can be a male name in an oblique case or a female name in the nominative. As oblique case nominals are a distinct and common feature of the Icelandic language (Thráinsson, 2007), this problem is highly relevant in terms of correctly relaying the meaning of the sentence (see Table 2).

## 4 Results

All submissions were scored using automatic metrics we constructed. Furthermore, we manually

reviewed around 150 randomly selected examples in the case of the idioms (around 100 'idiomatic' examples and around 50 'literal' examples for each submitted system). The authors reviewed the translations themselves, manually changing the scores given by our automatic method (using the 'positive' and 'negative' keywords discussed in 3.1) if they deemed it wrong.

The translations of our proper names suite was only carried out with naive automatic methods. The translations were lemmatized using a lemmatizer for Icelandic (Ingólfsdóttir et al., 2019) and compared with a reference of which Icelandic lemmas should appear in the translation and in which grammatical form (being able to look up lemmas is especially useful for the given names, since the male and female names share surface forms).

We show the results of our manual evaluation in Table 3 and the results of automatic metrics for our idioms test suite in Table 4. For our names test suite, we show the results of our automatic metrics in Table 5. Our scripts for running the automatic evaluations and the manually reviewed examples are released along with our test sets.

### 4.1 Scores for Idiomatic Expressions

Our results show a wide range of performance across different models. The best overall accuracy on the idioms test suite is achieved by Claude 3.5, with Unbabel-Tower70B a close second, as indicated both by our automatic and manual evaluation. Claude 3.5 is also the highest-scoring submission when we only consider translations of expressions used in an idiomatic sense, both according to our automatic metrics and the manual review, and Unbabel-Tower70B the clear runner-up.

When considering the literal translations in isolation, however, the overall two best models are narrowly 'beaten' by a few models that score considerably lower overall. According to our automatic

| System | Total Idioms | Total Literals | Idiom Accuracy | Literal Accuracy | Total Accuracy |
|---|---|---|---|---|---|
| AMI | 100 | 65 | 0.29 | 0.892308 | 0.527273 |
| Aya23 | 93 | 49 | 0.0537634 | 0.122449 | 0.0774648 |
| Claude-3.5 | 96 | 56 | **0.75** | 0.857143 | **0.789474** |
| CommandR-plus | 93 | 47 | 0.0967742 | 0.382979 | 0.192857 |
| CycleL | 92 | 49 | 0 | 0.102041 | 0.035461 |
| Dubformer | 91 | 53 | 0.340659 | 0.603774 | 0.4375 |
| GPT-4 | 93 | 48 | 0.430108 | 0.833333 | 0.567376 |
| IKUN-C | 95 | 52 | 0.494737 | 0.75 | 0.585034 |
| IKUN | 95 | 51 | 0.526316 | 0.607843 | 0.554795 |
| IOL_Research | 92 | 47 | 0.434783 | 0.702128 | 0.52518 |
| Llama3-70B | 93 | 50 | 0.268817 | 0.62 | 0.391608 |
| ONLINE-A | 188 | 107 | 0.265957 | 0.859813 | 0.481356 |
| ONLINE-B | 102 | 69 | 0.22549 | 0.898551 | 0.497076 |
| ONLINE-G | 97 | 66 | 0.185567 | 0.80303 | 0.435583 |
| TranssionMT | 76 | 50 | 0.223684 | **0.92** | 0.5 |
| TSU-HITs | 92 | 48 | 0.0434783 | 0.104167 | 0.0642857 |
| Unbabel-Tower70B | 95 | 57 | 0.631579 | 0.877193 | 0.723684 |

Table 3: Results of manual evaluation of system performance on our idioms test suite. We randomly split up the translations of the test suite into segments of around 100 'idiomatic' example translations and around 50 'literal' example translations (see 'Total' columns). The highest scores in each column are in bold. The authors reviewed the translations themselves and the reviewed examples, along with our grading, can be found at `https://github.com/stofnun-arna-magnussonar/idioms_names_test_suite/idioms/human_evaluation`.

| System name | Total score | Correct idiomatics | CI ratio | Correct literals | CL ratio |
|---|---|---|---|---|---|
| AMI | 0.447236 | 83 | 0.21 | 184 | **0.9** |
| Aya23 | 0.169179 | 39 | 0.1 | 62 | 0.3 |
| Claude-3.5 | **0.654941** | 216 | **0.55** | 175 | 0.86 |
| CommandR-plus | 0.293132 | 66 | 0.17 | 109 | 0.53 |
| CycleL | 0.108878 | 22 | 0.06 | 43 | 0.21 |
| Dubformer | 0.427136 | 112 | 0.28 | 143 | 0.7 |
| GPT-4 | 0.547739 | 161 | 0.41 | 166 | 0.81 |
| IKUN-C | 0.480737 | 141 | 0.36 | 146 | 0.72 |
| IKUN | 0.509213 | 161 | 0.41 | 143 | 0.7 |
| IOL_Research | 0.482412 | 133 | 0.34 | 155 | 0.76 |
| Llama3-70B | 0.417085 | 99 | 0.25 | 150 | 0.74 |
| ONLINE-A | 0.442211 | 86 | 0.22 | 178 | 0.87 |
| ONLINE-B | 0.447236 | 85 | 0.22 | 182 | 0.89 |
| ONLINE-G | 0.413735 | 71 | 0.18 | 176 | 0.86 |
| TranssionMT | 0.448911 | 86 | 0.22 | 182 | 0.89 |
| TSU-HITs | 0.112228 | 24 | 0.06 | 43 | 0.21 |
| Unbabel-Tower70B | 0.60804 | 195 | 0.5 | 168 | 0.82 |

Table 4: Results of automatic evaluation of system performance on our idioms test suite. We show the overall score for each system but also consider separately the percentage of idiomatic text examples marked as correct and the percentage of literals marked correct, to try and give an overview of the relationship between the two. Highest scores in each column are in bold. Our scripts for running automatic evaluation can be found at `https://github.com/stofnun-arna-magnussonar/idioms_names_test_suite/idioms`.

| System name | Total score | Total city score | Total people score |
|---|---|---|---|
| AMI | **0.5399** | **0.4705** | 0.5861 |
| Aya23 | 0.3838 | 0.0432 | 0.6103 |
| Claude-3.5 | 0.5091 | 0.4591 | 0.5423 |
| CommandR-plus | 0.3339 | 0.1205 | 0.4758 |
| CycleL | 0.0 | 0.0 | 0.0 |
| Dubformer | 0.4383 | 0.3614 | 0.4894 |
| GPT-4 | 0.5109 | 0.2773 | **0.6662** |
| IKUN-C | 0.4691 | 0.2727 | 0.5997 |
| IKUN | 0.4846 | 0.2886 | 0.6148 |
| IOL_Research | 0.4773 | 0.2205 | 0.648 |
| Llama3-70B | 0.4138 | 0.3227 | 0.4743 |
| ONLINE-A | 0.5345 | 0.4659 | 0.5801 |
| ONLINE-B | 0.5109 | 0.4273 | 0.5665 |
| ONLINE-G | 0.4065 | 0.3614 | 0.4366 |
| TranssionMT | 0.5082 | 0.4227 | 0.565 |
| TSU-HITs | 0.147 | 0.0932 | 0.1828 |
| Unbabel-Tower70B | 0.5254 | 0.4114 | 0.6012 |

Table 5: Results of automatic evaluation of system performance on our names test suite, given as a proportion of properly scored city names 'Total city score', properly scored given names 'Total people score' and overall 'Total score'. Highest scores in each column are in bold. Our scripts for running automatic evaluation can be found at https://github.com/stofnun-arna-magnussonar/idioms_names_test_suite/names. (Note that the zeroes for CycleL's submission are not a mistake, this submission performed poorly and our scoring strategy is not particularly forgiving.)

metrics, our own submission (AMI) scores highest in that category, only slightly ahead of ONLINE-B and TranssionMT. These three also come out on top in the manual evaluation, with TranssionMT recording the highest score (a superb 0.92) and ONLINE-B and AMI following in second and third.

This discrepancy between performance in translating phrases in an idiomatic context and a literal context is very interesting - these three models all scored under 0.3 in idiomatic accuracy, which suggests that for some models, proficiency in effectively translating text in a literal sense comes at a cost to their ability to handle more metaphorical text. The best-performing models overall, however, were seemingly able to maneuver quite effectively between both use cases. Models, perhaps predictably, generally score higher when translating literal usage than when translating idioms.

## 4.2 Scores for Proper Names

In terms of the proper names suite, place names prove to be much more difficult for the submitted models than people's names. It is the submission by our own team which narrowly tops the list overall, ahead of ONLINE-A and Unbabel. The AMI submission also ranks highest when place names are considered in isolation, although it still gets fewer than half of all names correct. For given names, GPT-4 scores highest.

For this part of our test set, we report no manual evaluation. A cursory glance at the output, however, shows that our naive automatic scoring method still leaves quite a bit to be desired. A problem with testing for specific grammatical forms in each case is that the correct form can change depending on the sentence structure. As discussed in 3.2, we tried to control for this by keeping test sentences brief and unambiguous. Even so, we find there are examples of different phrasings than we expected in some translation outputs that call for a different grammatical form of a name than our scoring mechanism supposes, but can still be considered a decent translation.

This especially applies to the sentence form: X "cares for" Y. We assumed a correct translation into Icelandic would be: X "þykir vænt um" Y, where X would take the dative case and Y the accusative. The submitted systems, however, had many different ideas on how best to phrase this system, not all of them completely wrong.

We therefore recognize that our scoring system needs to be fine-tuned but nevertheless believe the very low scores are mainly a reflection of the difficulty of this task.

## 5 Conclusions and Future Work

Scores on both sets are relatively low, indicating that these particular categories continue to pose some problems for even state-of-the-art translation models and that there is considerable room for improvement.

Future work can explore further comparison of performance and fine-tuning of our automatic scoring methods. Given time, we could also have investigated whether more manual evaluation, ideally using more evaluators, would have resulted in different scores.

We also note that our test suite can be adapted with relative ease into other languages and hope that this allows for further work on other language directions.

## Limitations

There are several judgment calls to be made when working with our chosen categories and many of the decisions we made in terms of selecting items to be translated, defining automatic metrics for 'right' and 'wrong' translations and manual evaluation can be argued for or against. We are aware that the choices we make could be indicative of potential biases of the authors and that a different team, perhaps with a different demographic makeup, might well have constructed the test set and evaluated the translations in a different way.

These necessary choices are perhaps most apparent in terms of our idioms set. Evaluation of linguistic acceptability of translations and correspondence of idiomatic phrases between languages is based on our intuition and we are aware that fluent speakers of English and Icelandic may disagree on some decisions. Another point to consider is the degree to which we want our test set to be prescriptive - as a simple search on the Internet can prove, there are multiple usages of common English idioms directly translated into Icelandic, e.g. on social media (Hilmisdóttir et al., 2023). Determining at what point to say this usage is no longer 'incorrect' is an interesting question of ethics and philosophy of language.

As for our set of proper names, there exists some speaker variation in how and when place names are translated into Icelandic, although we have tried to limit our set to fairly uncontroversial choices (see discussion in 3.2). The requirement of not translating English names into Icelandic is less cut and dried, as it may be appropriate for a machine translation model in some cases, e.g. in literary text or the discussion of royal or historical figures. It can also be noted that some of our English names are, in fact, given names in Iceland. This should not affect our results, however, as we allow for the inflection of a final -"a" into -"u" in female names like "Pamela" and in other cases, 'non-Icelandic' names typically remain completely unchanged in all grammatical cases.

## References

Tosin Adewumi, Roshanak Vadoodi, Aparajita Tripathy, Konstantina Nikolaido, Foteini Liwicki, and Marcus Liwicki. 2022. Potential idiomatic expression (PIE)-English: Corpus for classes of idioms. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 689–696, Marseille, France. European Language Resources Association.

Kristín Bjarnadóttir, Kristín Ingibjörg Hlynsdóttir, and Steinþór Steingrímsson. 2019. DIM: The Database of Icelandic Morphology. In *Proceedings of the 22nd Nordic Conference on Computational Linguistics*, Turku, Finland.

Marzieh Fadaee, Arianna Bisazza, and Christof Monz. 2018. Examining the tip of the iceberg: A data set for idiom translation. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Hessel Haagsma, Johan Bos, and Malvina Nissim. 2020. MAGPIE: A large corpus of potentially idiomatic expressions. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 279–287, Marseille, France. European Language Resources Association.

Björn Halldórsson, Árni Davíð Magnússon, Finnur Ágúst Ingimundarson, Einar Freyr Sigurðsson, Steinþór Steingrímsson, Halldóra Jónsdóttir, and Þórdís Úlfarsdóttir. 2022. Idiomatic expressions (Icelandic and English) 22.09. CLARIN-IS.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita,

Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? A comprehensive evaluation.

Helga Hilmisdóttir, Martina Huhtamäki, and Susanna Karlsson. 2023. Pragmatic borrowing from English. *Nordic Journal of Linguistics*, 46(3):255–256.

Svanhvít Lilja Ingólfsdóttir, Hrafn Loftsson, Jón Friðrik Daðason, and Kristín Bjarnadóttir. 2019. Nefnir: A high accuracy lemmatizer for Icelandic. In *Proceedings of the 22nd Nordic Conference on Computational Linguistics*, pages 310–315, Turku, Finland. Linköping University Electronic Press.

Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.

Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, He Wang, Renlong Ai, Shushen Manakhimova, Ursula Strohriegel, Sebastian Möller, and Hans Uszkoreit. 2022. A linguistically motivated test suite to semi-automatically evaluate German–English machine translation output. In *Proceedings of the Language Resources and Evaluation Conference*, pages 936–947, Marseille, France. European Language Resources Association.

Vikas Raunak, Arul Menezes, Matt Post, and Hany Hassan. 2023. Do GPTs produce less literal translations? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1041–1050, Toronto, Canada. Association for Computational Linguistics.

Vered Shwartz. 2021. A long hard look at MWEs in the age of language models. In *Proceedings of the 17th Workshop on Multiword Expressions (MWE 2021)*, page 1, Online. Association for Computational Linguistics.

David Stap, Eva Hasler, Bill Byrne, Christof Monz, and Ke Tran. 2024. The fine-tuning paradox: Boosting translation quality without sacrificing LLM abilities. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6189–6206, Bangkok, Thailand. Association for Computational Linguistics.

Steinþór Steingrímsson, Einar Freyr Sigurðsson, and Björn Halldórsson. 2024. Evaluating Capabilities of MT Systems in Translating Idiomatic Expressions Using a Specialized Dataset. In *Proceedings of CLARIN annual conference 2024, October 15-17, 2024*.

Kenan Tang. 2022. Petci: A parallel English translation dataset of Chinese idioms.

Höskuldur Thráinsson. 2007. *The Syntax of Icelandic*. Cambridge University Press, Cambridge.

Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2024. A paradigm shift in machine translation: Boosting translation performance of large language models. In *The Twelfth International Conference on Learning Representations*.