# Scrambling in Persian from English Subtitles

John R. Starr

# Presentation Outline:

1. Quick Description of Persian and Scrambling
2. Information on the Tehran English-Persian Parallel Corpus
3. Reorganization of Data
4. Modification of Data
5. Analysis of Data
6. Conclusions

# Part I:
# On Persian

General characteristics, example
sentence, scrambling, project motivation

# Persian: A Brief Description
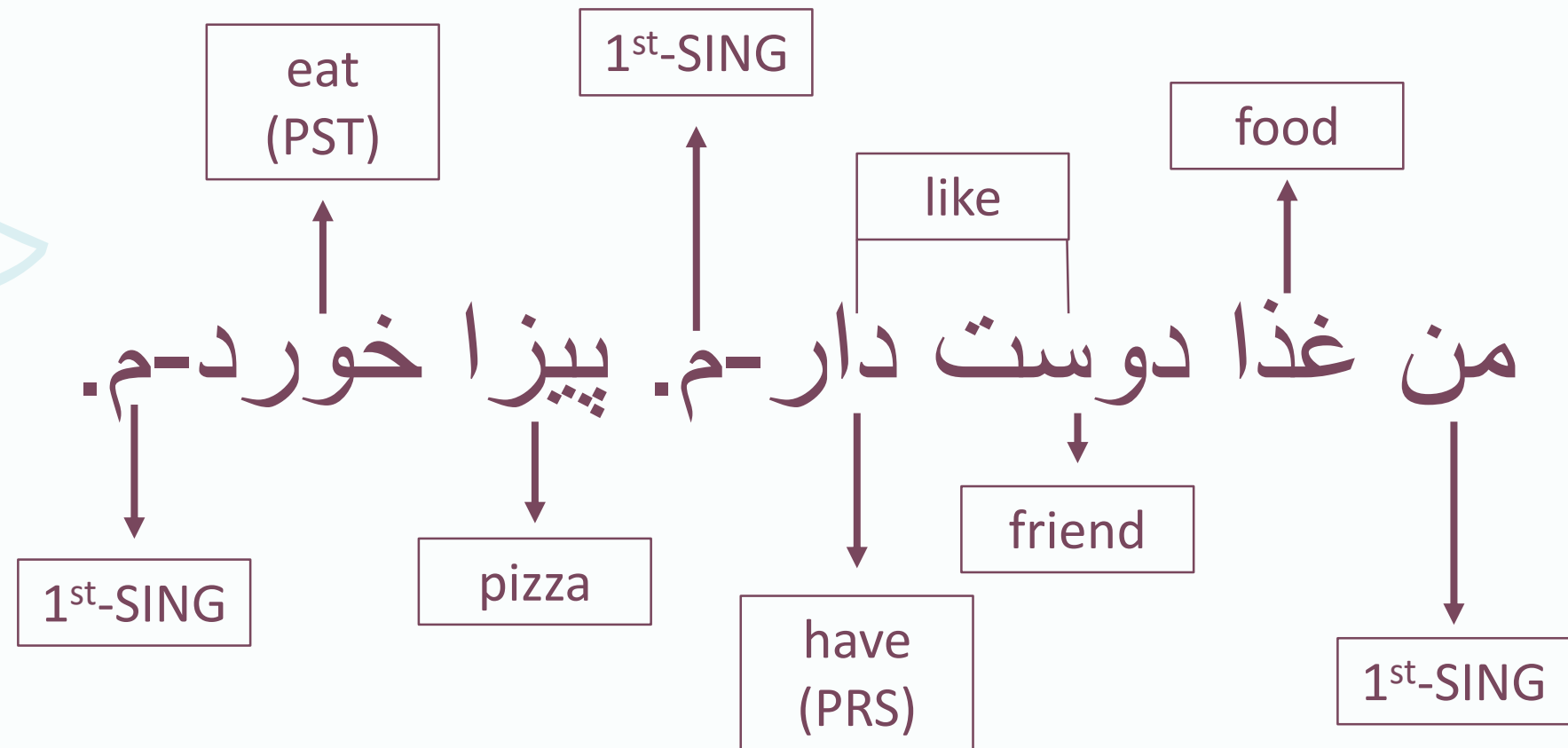
– Indo-Iranian language spoken in Iran, Afghanistan, and Tajikistan

– Language subfamilies: Farsi (this project), Dari, and Tajik

– Reads right-to-left, using the Arabic alphabet (+4 letters)

– Many verbs are compound and include non-verbal components

– Is "pro-drop": subject pronouns optional if already stated

– Underlyingly SOV, but is open to scrambling (DEFINE SCRAMBLING)

# An Example Sentence:

eat
(PST)

1<sup>st</sup>-SING

food

like

من غذا دوست دار-م. پیزا خورد-م.

1<sup>st</sup>-SING

pizza

have
(PRS)

friend

1<sup>st</sup>-SING

a. weil mich die anderen oft einladen

b. weil mich die anderen mich oft einladen
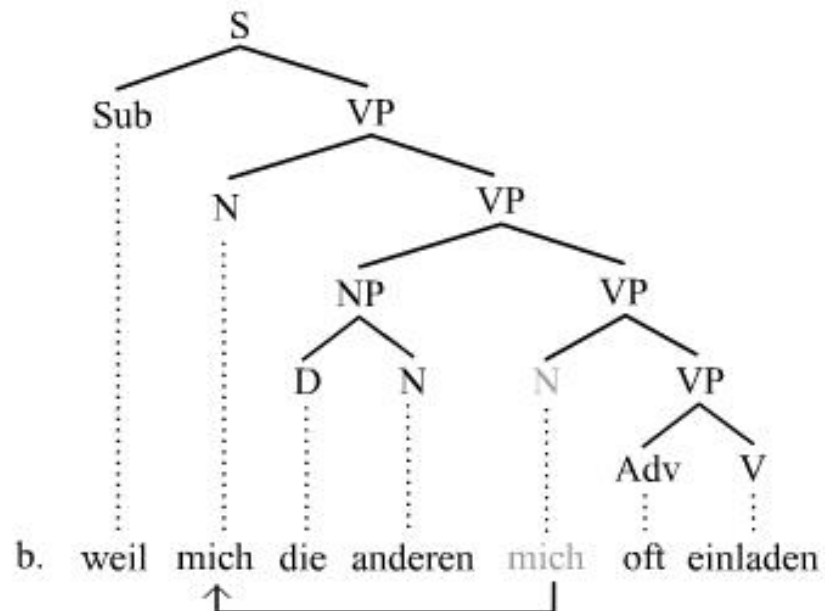
# What is scrambling?

An example from German: "Because the others invite me"

# Part II:
# The Tehran English-Persian Parallel Corpus

Layout of corpus, example of raw data, positives/negatives of corpus

# The Tehran English-Persian Parallel Corpus: What's it all about?

– Parallel corpus of English-Persian subtitles

– 554, 621 lines total

 – About 3.8 million words for each language

– Obtained from approx. 21,000 movie files from Open-subtitles (free online collection of movie subtitles in various languages)

– One large .xml file for each language:

The Raw Data

```xml
<?xml version="1.0" encoding="utf-8"?>
<letsmt version="1.0">
<head></head>
<body>
<s id="1">raspy breathing .</s>
<s id="2">dad .</s>
<s id="3">maybe its the wind .</s>
<s id="4">no .</s>
<s id="5">stop please stop .</s>
<s id="6">you have a week , evans then well burn the house .</s>
<s id="7">william .</s>
<s id="8">god damn it , william .</s>
<s id="9">god damn it put that down .</s>
<s id="10">let go .</s>
<s id="11">its the last feed weve got .</s>
```

# The Ups and Downs of the TEP

## +

– Data intended to sound like regular speech, rather than formal writing

– Accurate alignment of speakers due to timestamps

– .xml files clearly indexed

## -

– Lack of punctuation

– Spoken Persian behaves very differently than written Persian

– Persian spelling varies greatly (but can have a great effect)

– Many incomplete sentences

# Part III: Organizing the Data

Extracting data from .xml files, moving to DFs, combining data into one big DF

# The Process:

```
In [16]:  full_df = pd.Series(eng_lines).to_frame('Eng').join(pd.Series(far_lines).to_frame('Far'), how='outer')
          full_df.index.name = 'ID'
          full_df
```

Out[16]:

| ID | Eng | Far |
|---|---|---|
| 1 | raspy breathing | صداي خر خر |
| 2 | dad | پدر |
| 3 | maybe its the wind | شايد صداي باد باشه |
| 4 | no | نه |
| 5 | stop please stop | دست نگه داريد خواهش ميکنم دست نگه داريد |

Fill & Line

```
In [1]:  1  import pandas as pd
         2  import numpy as np
         3  import nltk
         4  # import postagger
         5  import xml.etree.ElementTree as ET
         6  from lxml import etree
```

```
In [6]:  1  eng_lines_test = {}
         2  for item in root_eng.findall('./body/s')[:5]:
         3      eng_lines_test[int(str(item.values()).replac
```

```
In [7]:  1  eng_lines_test.keys()
```

Out[7]: dict_keys([1, 2, 3, 4, 5])

```
In [8]:  1  eng_lines_test.values()
```

Out[8]: dict_values(['raspy breathing', 'dad', 'maybe its the

```
In [21]:  full_df.head()
```

Out[21]:

| ID | Eng | Far | Eng_Tok | Far_Tok | Eng_Len | Far_Len | Eng_Types | Far_Types |
|---|---|---|---|---|---|---|---|---|
| 1 | raspy breathing | صداي خر خر | [raspy, breathing] | [صداي, خر, خر] | 2 | 3 | {breathing, raspy} | {صداي, خر} |
| 2 | dad | پدر | [dad] | [پدر] | 1 | 1 | {dad} | {پدر} |
| 3 | maybe its the wind | شايد صداي باد باشه | [maybe, its, the, wind] | [شايد, صداي, باد, باشه] | 4 | 4 | {wind, the, maybe, its} | {شايد, باشه, صداي, باد} |
| 4 | no | نه | [no] | [نه] | 1 | 1 | {no} | {نه} |

Fill & Line

```
In [25]:  full_df.Eng_Len.value_counts()
```

```
Out[25]:  6    60403
          1    59483
          7    58117
          5    57950
          8    54154
          4    52814
          9    47822
          10   40727
          3    40255
          2    36142
```

```
In [26]:  full_df.Far_Len.value_counts()
```

```
Out[26]:  2    67818
          5    60178
          6    60088
          4    56988
          3    55920
          7    55866
          8    49803
          9    42446
          1    35506
          10   34563
```

# Part IV:
# Modifying the Data

---

POS-tagging, chunking, generalizing

# POS-tagging and Chunking Persian Text

```
In [0]:   # Importing the necessary modules
          from hazm import *
```

```
In [3]:   # Building our tagger
          tagger = POSTagger(model='postagger.model')
          tagger.tag(word_tokenize('ما بسیار کتاب می‌خوانیم'))
```

```
Out[3]:   [('ما', 'PRO'), ('بسیار', 'ADV'), ('کتاب', 'N'), ('می\u200cخوانیم', 'V')]
```

```
In [0]:   # Building our tagger
          chunker = Chunker(model='chunker.model')
```

```
In [5]:   # Test file for the chunker
          tagged = tagger.tag(word_tokenize('کتاب خواندن را دوست داریم'))
          tree2brackets(chunker.parse(tagged))
```

```
Out[5]:   '[خواندن کتاب NP] [را POSTP] [دوست داریم VP]'
```

# How to properly generalize the WO?

– Data is in strings → Determine WO with regex

  – Alternative methods?

– Defining NPs/VPs <u>and</u> defining what they're *not*:

  – What is:

```
In [21]:   1  np = r'\[[^\]]+ NP\]'
           2  vp = r'\[[^\]]+ VP\]'
```

  – What isn't:

```
In [22]:   1  np_not = r'\[[^\]]+ [^N]P\]'
           2  vp_not = r'\[[^\]]+ [^V]P\]'
```
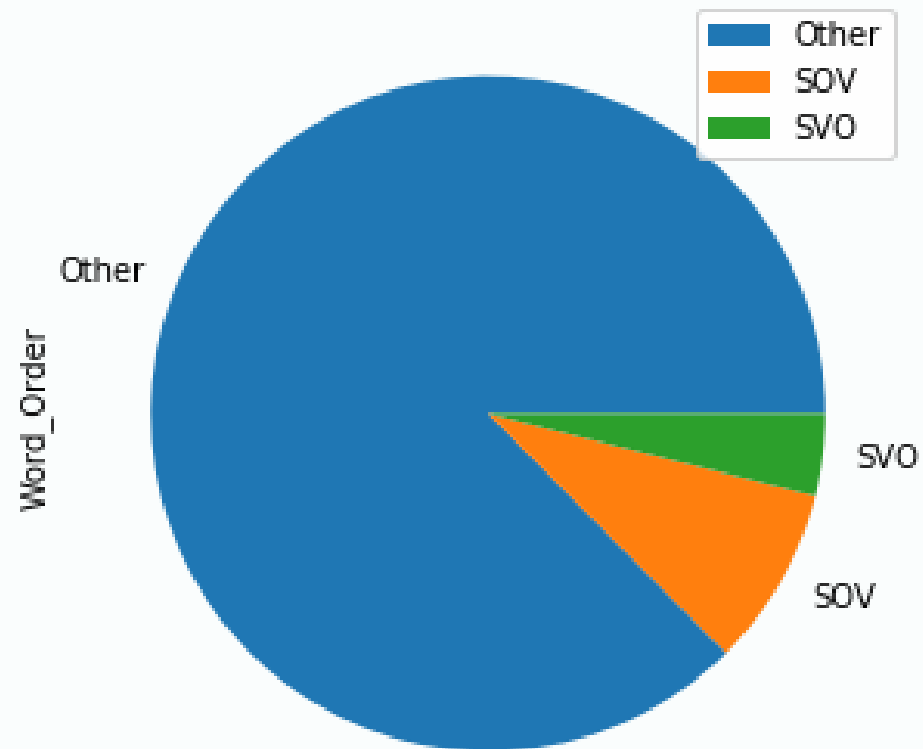
# Part V:
# Analyzing the Data

How the data is split up, general remarks,
looking at data using the word "like"

How do the
data look?

Yikes…

# Some Qualities of Generalizations:

– The chunker is… not great.

– SOV/SVO sentences generally correct (counts not yet discovered)

– "Other sentences" are either too simple, too complex, or have a weird quirk about them

– Looking into "is like" sentences:

# Example SOV Sentences:

```
In [36]:  for item in like_sov_text[:10]:
              print(item)
              print()
              print()
```

[است VP] [NP استاره 5] [NP یک هتل] [PP مثل] [NP قلعه کاپریکورن]

[NP میری؟] [PP داري] [ADVP کجا] [VP کني] [NP فکر مي] [PP ،باز] [NP هي آدم حقه]

[NP این دنیا مهربوني میبخشه] [PP به] که [VP الهیه] [NP یک فرشته] [PP مثل] [NP اون زن]

[NP رنگ طلاست] [PP به] که [VP خورشید] [NP نور] [PP براي] [For sunlight is like gold NP]

[VP توش داره] [ADVP زیادي] [VP رازهاي] [NP جرم] [NP صحنه]

[NP شنه] [VP هاي] [NP دانه] [PP مثل] [NP قدرت]

# Example SVO Sentences:

```
In [38]:    1  for item in like_svo_text[:10]:
            2      print(item)
            3      print()
            4      print()
```

[NP اغاز یه مسابقه تلوزیونیه] [VP شبیه] [NP این ویزیتها]

[NP مثله خوره میمونه] [VP آهنگه] [NP این]

[NP هنر کفش] [VP هست] [NP این]

[NP خلاصه میشه] [NP بابا] و [NP مامان] [VP بین] [NP کورس من] [VP توي] [NP کریسمس]

[NP یک گل رز قرمزه] [VP شبیه] [NP عشقمون]

[VP کني] [NP زیاد اشتباه] [VP نباید] [VP نباید] [NP فضاست] [NP زندگي تو] [VP شبیه] [NP زندگي این پایین]

[NP پایمال کنیم] [POSTP را] [NP درست شبیه اینه که خون شهدا] [VP شرم آور] [VP هاي] [NP رفتار]

# What about the "Other" sentences?

```
In [47]:   1  for item in other_sample.sample(10):
           2      print(item)
           3      print()
```

[NP وقتي كه  اشراف جلسه تشكيل بدن] [PP تا]

[NP چارلي] [VP ممنونم] .

[VP ت] [POSTP را] [NP دارم] [POSTP را] [VP ت] [ADVP هم] [NP من] [VP رفتم] [NP اون خونه] [PP به] [NP من] كه [NP وقتی] [PP از] [NP یک همچین چیزایی]
[NP جربه میکنم

[NP اون پسر خشمگین شدن] [PP نسبت به] [NP اونها]

[NP تخم هارو داغ و خوب ميكنيم اينجوري سفت و محكم ميشن] [NP ما] [ADVP حالا]

[NP ميكنن انجام ميدي] [PP بقيه] [NP هر كاري كه] [NP تو] [VP نداشتي] [NP چند ماهه وضع خوبي] [NP تو] [ADVP پس]

و [NP من] [NP اونو] [VP اين شكلي] ميبينم

[VP اون اومد] ، [VP باش] [NP هان آروم] ، [NP توليور] ، [NP بين شما] [VP بوول بياد] [VP گذاشتيد] [NP چه علت] [PP به]
NP]

[VP مى شناسند؟] [POSTP را] [NP همديگر] [NP آنجا] [PP از] [NP مينگ جو] و [NP يانگوم] [ADVP پس]

[NP غمزده عشق]

# Part VI: Conclusions

---

(Preliminary) results, further directions,

possible corrections

# Preliminary Results

– Constructions from similes lean towards SOV, whereas metaphorical constructions lean towards SVO

– More abstract comparisons also lean towards SVO

– Many different constructions will be analyzed

# Limitations of Project:

– Persian's inconsistent transcription and subject availability

– Persian parser/chunker/generalizer are limited and do not appear to have high levels of accuracy

–  Nearly 80% of the data is categorized as "Other"

# Further Directions:

- Machine learning program that categorizes sentences for scrambling

- Working with the speech data itself

- Sentiment marking

# THE END

Questions?

# References:

– Link to Tehran-Parallel Corpus: http://opus.nlpl.eu/TEP.php

  – Information on the TEP: https://link.springer.com/content/pdf/10.1007%2F978-3-642-19437-5.pdf

  – M. T. Pilevar, H. Faili, and A. H. Pilevar, â€œTEP: Tehran English-Persian Parallel Corpusâ€• in proceedings of 12th International Conference on Intelligent Text Processing and Computational Linguistics (CICLing-2011).

– Information on Hazm: https://github.com/sobhe/hazm

– Project Github: https://github.com/Data-Science-for-Linguists-2019/Scrambling-in-English-to-Persian-Subtitles