

# Still all Greekl-ish to me: Greekl-ish to Greek Transliteration

Anastasios Toumazatos, John Pavlopoulos,  
Ion Androutsopoulos, Stavros Vassos



# Greeklish



- Greek written in the Latin alphabet.
- Very prevalent in online informal settings.

# Greeklish



- Greek written in the Latin alphabet.
- Very prevalent in online informal settings.
- No standardized mapping between Greeklish and Greek.

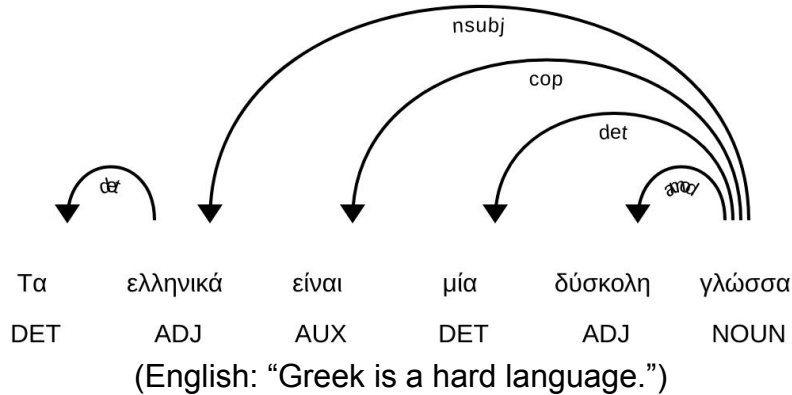
| Greeklish | Greek | Category    |
|-----------|-------|-------------|
| e         | αι    | phonetic    |
| e         | ε     | key-sharing |
| ai        | αι    | optical     |
| x         | χ     | optical     |
| x         | ξ     | phonetic    |
| ch        | χ     | phonetic    |
| th        | θ     | phonetic    |
| u         | θ     | key-sharing |
| 8         | θ     | optical     |

# Greeklish

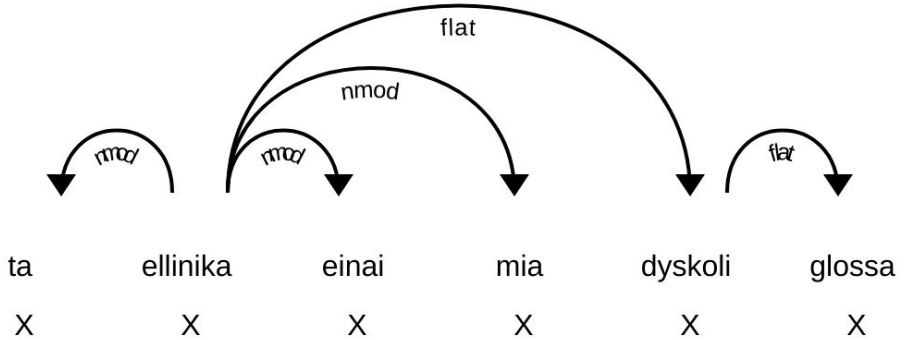
- Lack of standardization makes Greeklish hard to handle in various applications.

# Greeklish

- Lack of standardization makes Greeklish hard to handle in various applications.
- Example: Dependency parsing with **spaCy** (same sentence in Greeklish and Greek).




(a) Greek



(a) Greeklish

# Data (synthetic)

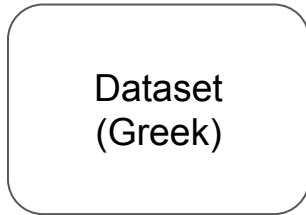
- No existing datasets for Greeklish-to-Greek (G2G) transliteration.
- We created new synthetic datasets, by converting Greek text data to Greeklish using a mapping of all plausible transliterations of Greek letters and choosing randomly each time.
- Data from **Europarl (1996-2011)** .
- Subtitles from two TV Shows, **Para 5** and **Friends (1st season)**.



Dataset  
(Greek)

# Data (synthetic)

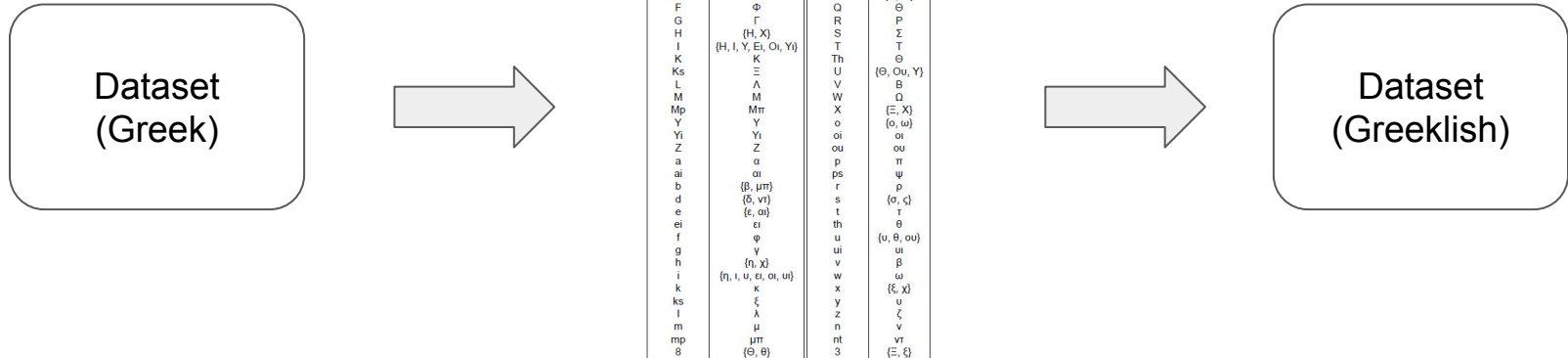
- No existing datasets for Greeklish-to-Greek (G2G) transliteration.
- We created new synthetic datasets, by converting Greek text data to Greeklish using a mapping of all plausible transliterations of Greek letters and choosing randomly each time.
- Data from **Europarl (1996-2011)** .
- Subtitles from two TV Shows, **Para 5** and **Friends (1st season)**.



| Greeklish | Greek              | Greeklish | Greek      |
|-----------|--------------------|-----------|------------|
| A         | Α                  | N         | Ν          |
| Ai        | Αι                 | Nt        | Ντ         |
| B         | {Β, Μπ}            | O         | {Ο, Ω}     |
| D         | {Δ, Ντ}            | Oi        | Οι         |
| E         | {Ε, Αι}            | Ou        | Ου         |
| Ei        | Ει                 | P         | {Π, Φ}     |
| F         | Φ                  | Q         | Θ          |
| G         | Γ                  | R         | Ρ          |
| H         | {Η, Χ}             | S         | Σ          |
| I         | {Ι, Υ, Ει, Οι, Υι} | T         | Τ          |
| K         | Κ                  | Th        | Θ          |
| Ks        | Ξ                  | U         | {Θ, Ου, Υ} |
| L         | Λ                  | V         | Β          |
| M         | Μ                  | W         | Ω          |
| Mp        | Μπ                 | X         | {Ξ, Χ}     |
| Y         | Υ                  | o         | {ο, ω}     |
| Yi        | Υι                 | oi        | οι         |
| Z         | Ζ                  | ou        | ου         |
| a         | α                  | p         | π          |
| ai        | αι                 | ps        | ψ          |
| b         | {β, μπ}            | r         | ρ          |
| d         | {δ, ντ}            | s         | {σ, ς}     |
| e         | {ε, αι}            | t         | τ          |
| ei        | ει                 | th        | θ          |
| f         | φ                  | u         | {υ, θ, ου} |
| g         | γ                  | ui        | υι         |
| h         | {η, χ}             | v         | β          |
| i         | {ι, υ, ει, οι, υι} | w         | ω          |
| k         | κ                  | x         | {ξ, χ}     |
| ks        | ξ                  | y         | υ          |
| l         | λ                  | z         | ζ          |
| m         | μ                  | n         | ν          |
| mp        | μπ                 | nt        | ντ         |
| s         | {σ, ς}             | 3         | {Ξ, Χ}     |

# Data (synthetic)

- No existing datasets for Greeklish-to-Greek (G2G) transliteration.
- We created new synthetic datasets, by converting Greek text data to Greeklish using a mapping of all plausible transliterations of Greek letters and choosing randomly each time.
- Data from **Europarl (1996-2011)** .
- Subtitles from two TV Shows, **Para 5** and **Friends (1st season)**.





# Data (real-world)

- We also created smaller test datasets by manually transliterating real-world Greeklish content to Greek.
- Transcribed by MSc students trained for annotation.

| Name        | Description  |
|-------------|--|
| Survivorbot | Greeklish data from an online chatbot deployed during the run of a popular TV show in Greece ('Survivor'). |
| Gazzetta    | Greeklish data from the discussion forum of a popular sports website ('Gazzetta').                         |

# RBSLM

Shmera o kairos eine kalos.  
(English: "Today the weather is nice.")

# RBSLM

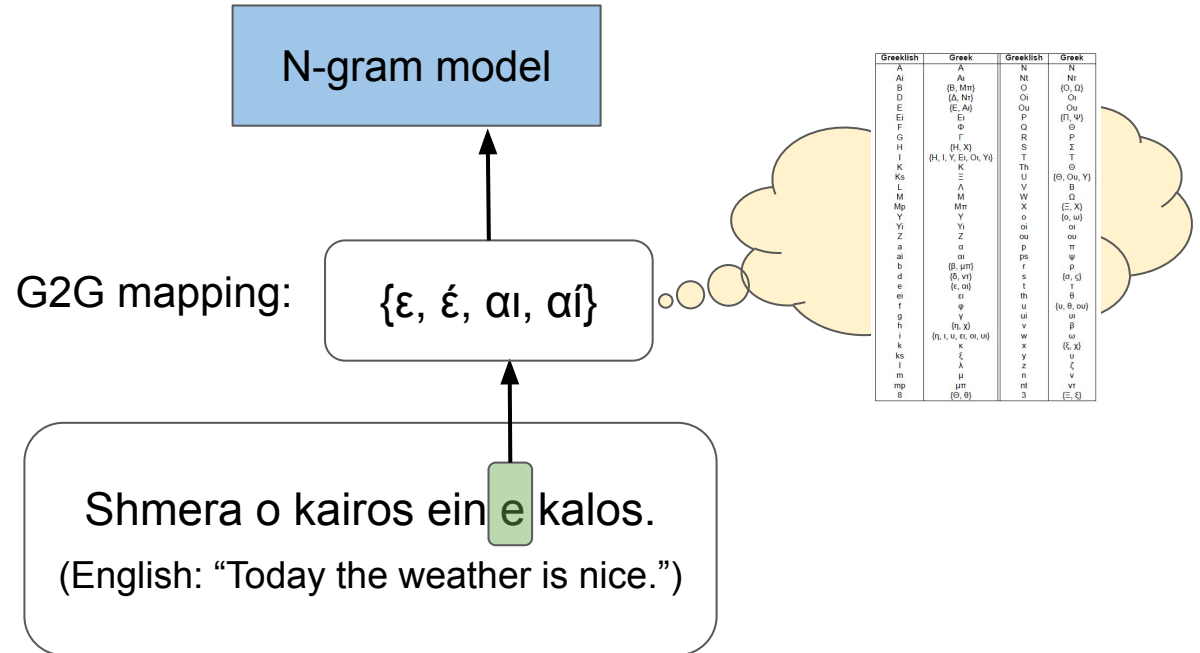
G2G mapping:

{ε, έ, αι, αί}

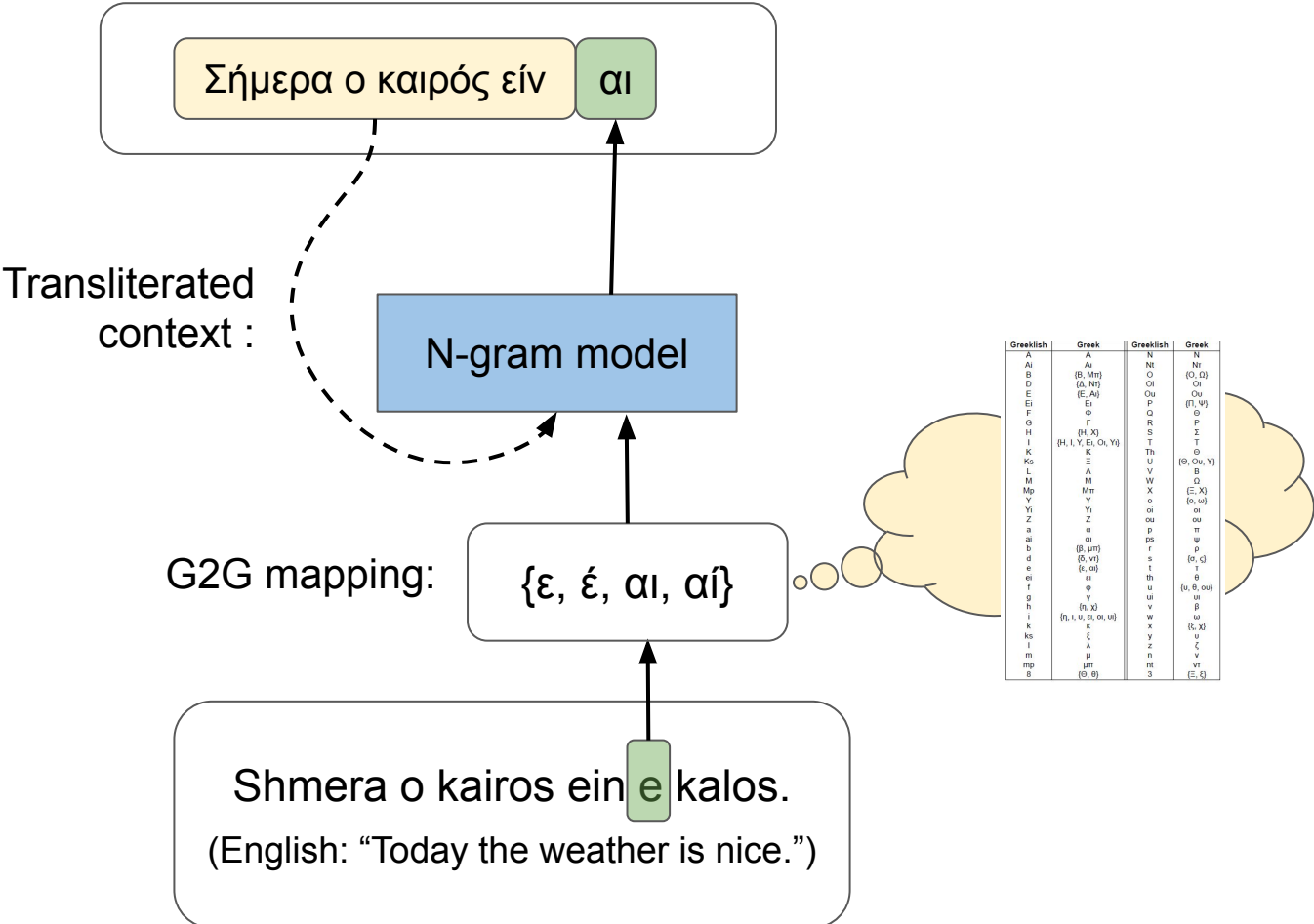
Shmera o kairos ein e kalos.  
(English: "Today the weather is nice.")

| Greeklish | Greek            | Greeklish | Greek      |
|-----------|------------------|-----------|------------|
| A         | Α                | N         | Ν          |
| Ai        | Αι               | Nt        | Ντ         |
| B         | (B, Mf)          | O         | (O, O)     |
| D         | (A, Nt)          | Oi        | Οι         |
| E         | (E, Ai)          | Ou        | Ου         |
| Ei        | Ει               | P         | (Pi, Pi)   |
| F         | φ                | Q         | Q          |
| G         | γ                | R         | Ρ          |
| H         | (H, X)           | S         | Σ          |
| I         | (I, Y, E, O, Yi) | T         | Τ          |
| K         | Κ                | Th        | Θ          |
| Ks        | Ξ                | U         | (O, O, Y)  |
| L         | Λ                | V         | Β          |
| M         | Μ                | W         | Ω          |
| Mp        | Mf               | X         | (Ξ, X)     |
| Y         | Υ                | o         | (o, o)     |
| Yi        | Υι               | oi        | οι         |
| a         | α                | ou        | ου         |
| ai        | αι               | pi        | πι         |
| b         | (B, Mf)          | ps        | ψ          |
| d         | (D, Vt)          | s         | (S, S)     |
| e         | (E, Ai)          | t         | τ          |
| ei        | ει               | th        | θ          |
| f         | φ                | u         | (u, O, ou) |
| g         | γ                | ui        | υι         |
| h         | (H, X)           | v         | β          |
| i         | (I, Y, E, O, Yi) | w         | ω          |
| ks        | Ξ                | x         | (X, X)     |
| k         | κ                | y         | Υ          |
| i         | ι                | z         | Ζ          |
| m         | μ                | v         | ν          |
| mp        | Mf               | vt        | ντ         |
| B         | (B, Mf)          | S         | (S, S)     |

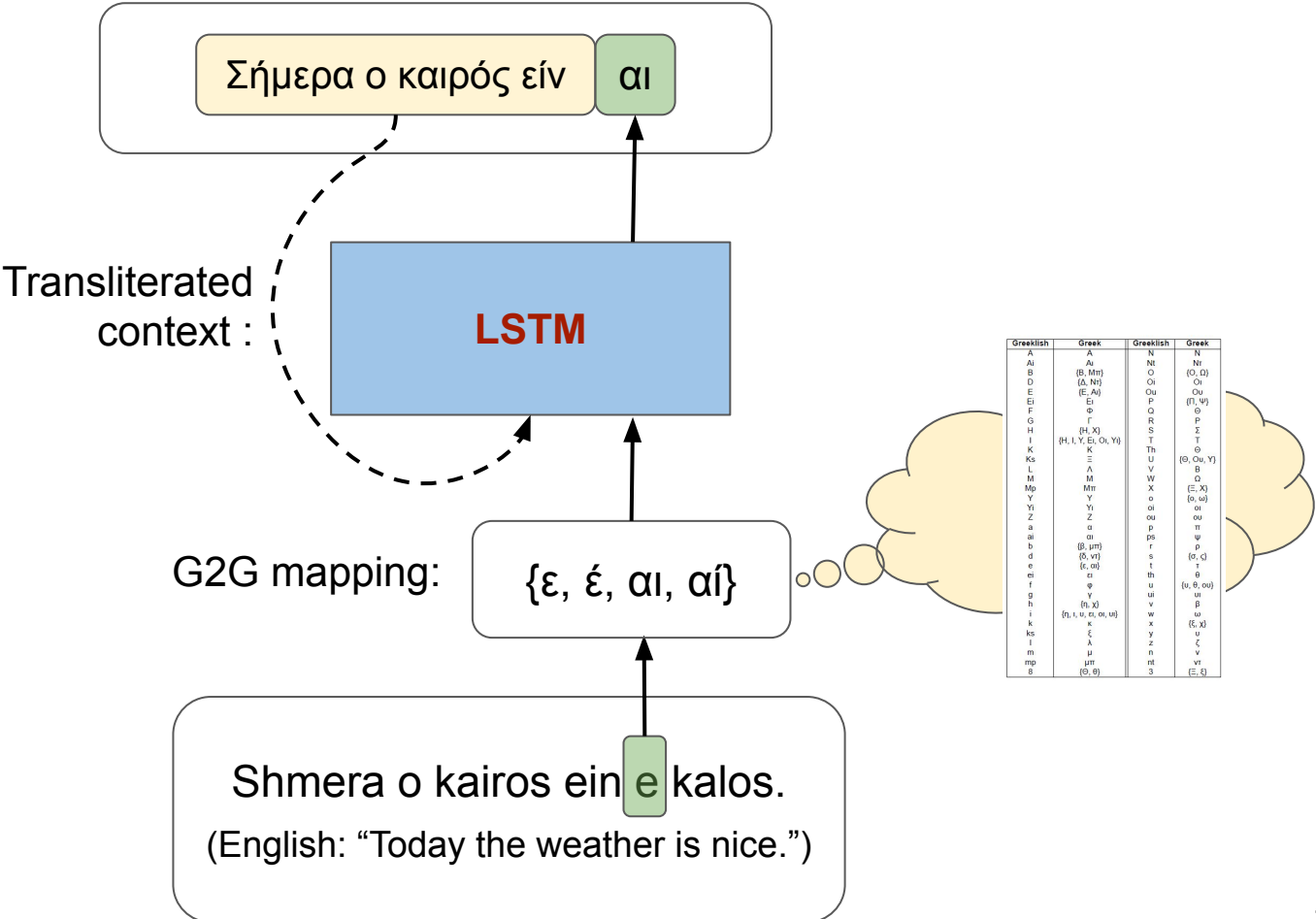
# RBSLM



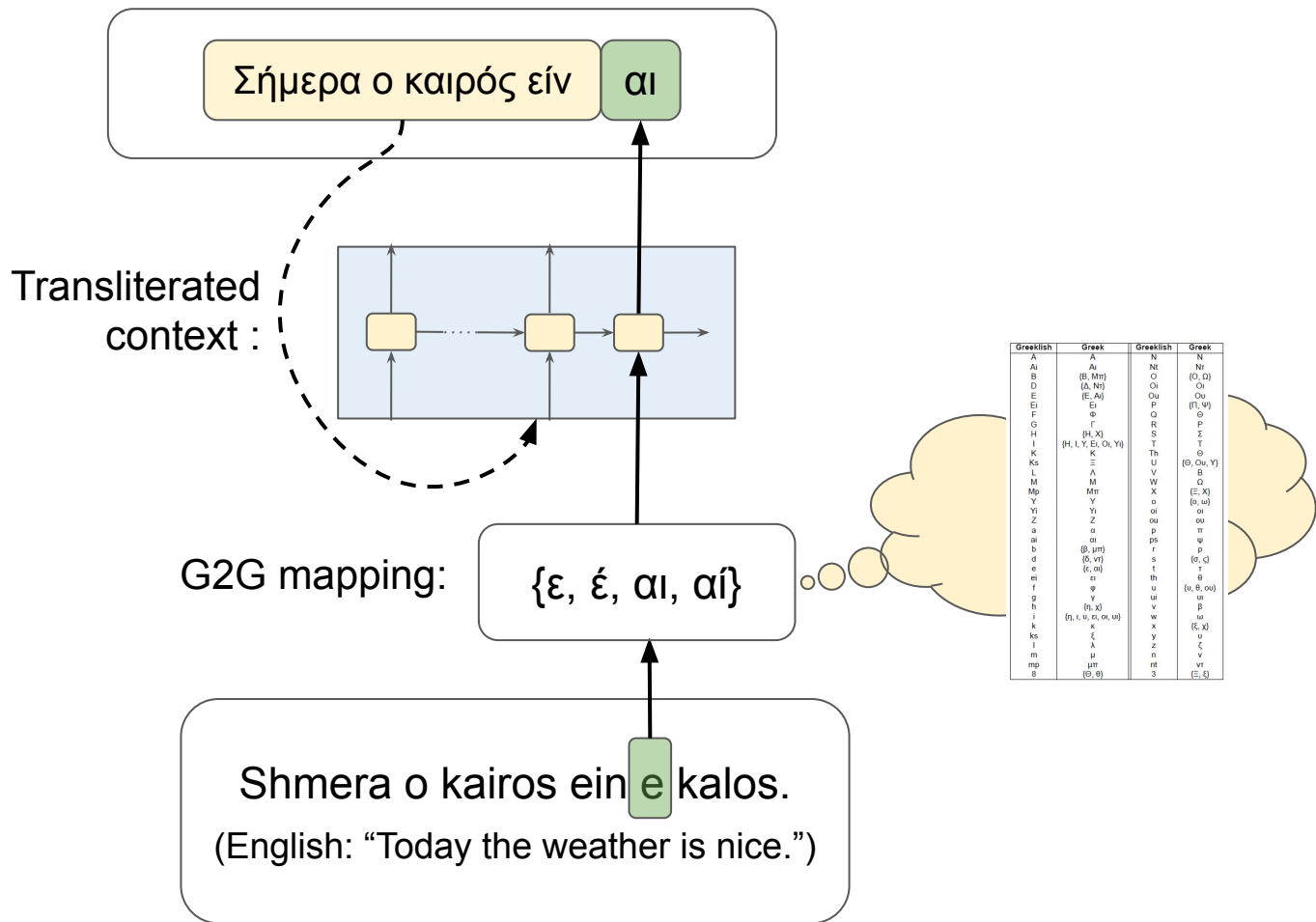
# RBSLM



# RBNLM



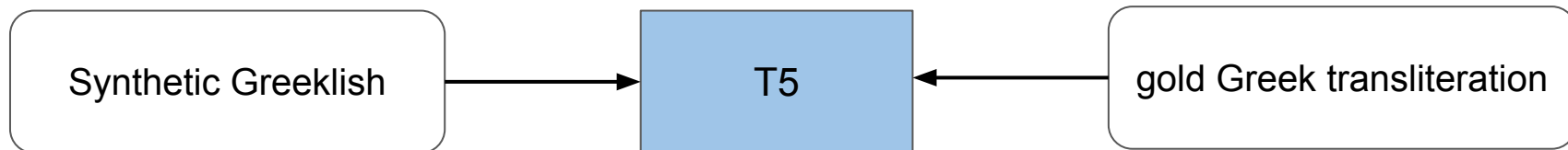
# RBNLM



# T5-based

- We used Google's **mT5** & **ByT5** models.
- Both are pre-trained on the same multilingual dataset.
- We further fine tune them on data from our *synthetic* datasets.

Fine-Tuning

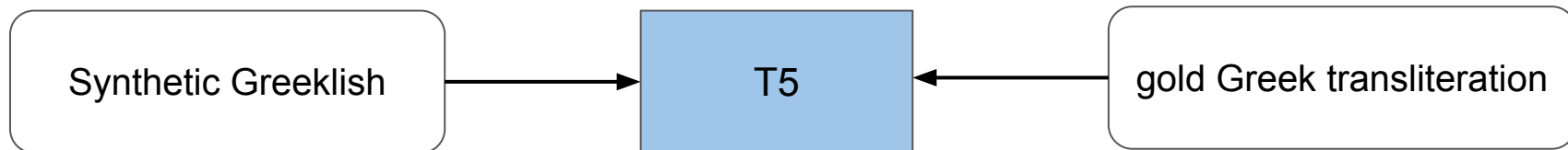




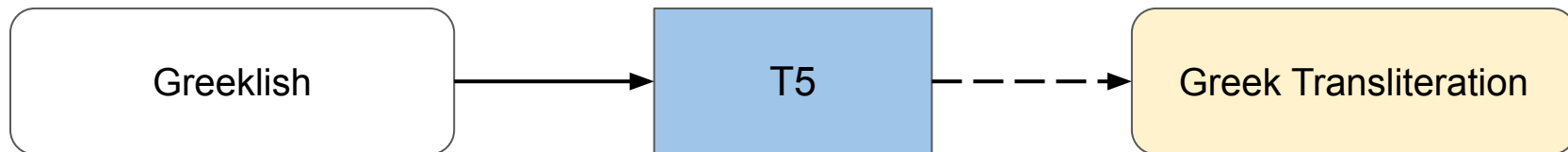
# T5-based

- We used Google's **mT5** & **ByT5** models.
- Both are pre-trained on the same multilingual dataset.
- We further fine tune them on data from our *synthetic* datasets.

Fine-Tuning



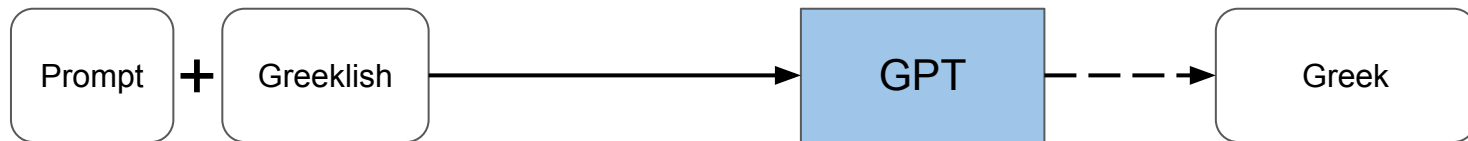
Usage



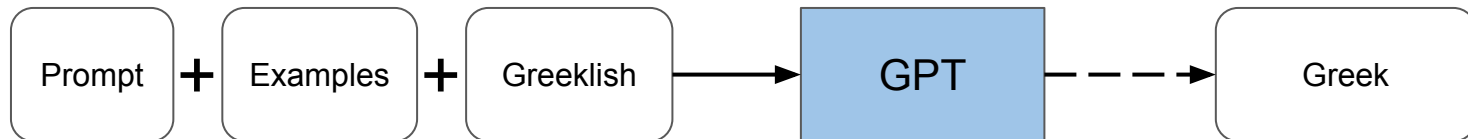
# GPT-based

- We used OpenAI's **GPT-3.5-t** and **GPT-4\***, also experimenting with the web version of ChatGPT.
- We investigated various prompts, as well as zero-shot and few-shot prompting. We find that few-shot prompting is helpful regardless of the relevance of topic (e.g. Europarl examples for transliterating a sports-related sentence).

Zero-shot

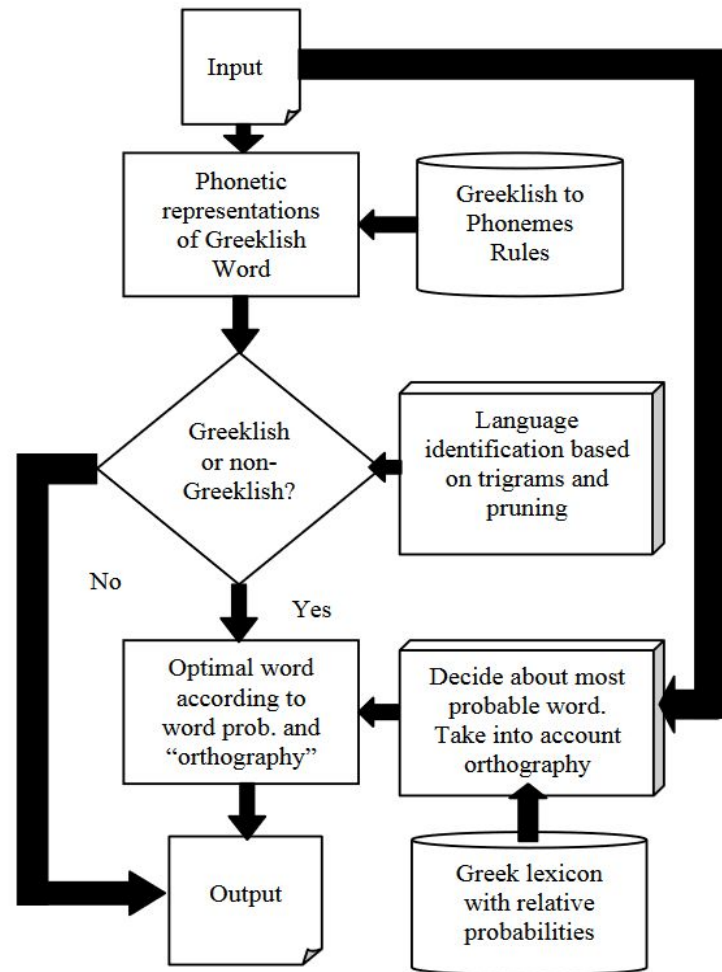


Few-shot



# All Greek to me!

- Proposed by Chalamandaris et al. (2006)
- Trained on real-world web-collected data.
- Based on a statistical tri-gram model & lexicons from large corpora.
- Figure of the architecture used from the original paper by Chalamandaris et al. (<https://aclanthology.org/L06-1229/>)



# Results (synthetic data)

| Model       | Europarl |       | Friends (TV) |       | Para5 (TV) |       |
|-------------|----------|-------|--------------|-------|------------|-------|
|             | CER      | WER   | CER          | WER   | CER        | WER   |
| AllGreek    | 2.98     | 11.54 | 7.56         | 19.62 | 5.32       | 16.32 |
| RBSLM       | 4.71     | 15.92 | 14.29        | 39.67 | 12.98      | 36.53 |
| RBNLM       | 1.34     | 4.45  | 9.68         | 27.03 | 7.36       | 20.55 |
| mT5         | 6.06     | 10.14 | 32.60        | 48.53 | 31.56      | 43.86 |
| ByT5        | 1.30     | 4.11  | 14.36        | 38.20 | 13.93      | 34.28 |
| ByT5-EU     | 0.64     | 2.30  | 11.36        | 29.74 | 9.80       | 25.04 |
| ByT5-TV     | 1.37     | 3.92  | 3.77         | 11.69 | 3.51       | 9.93  |
| GPT-4-6shot | 2.27     | 5.21  | 9.74         | 21.41 | 8.34       | 25.15 |

# Results (synthetic data)

| Model       | Europarl |       | Friends (TV) |       | Para5 (TV) |       |
|-------------|----------|-------|--------------|-------|------------|-------|
|             | CER      | WER   | CER          | WER   | CER        | WER   |
| AllGreek    | 2.98     | 11.54 | 7.56         | 19.62 | 5.32       | 16.32 |
| RBSLM       | 4.71     | 15.92 | 14.29        | 39.67 | 12.98      | 36.53 |
| RBNLM       | 1.34     | 4.45  | 9.68         | 27.03 | 7.36       | 20.55 |
| mT5         | 6.06     | 10.14 | 32.60        | 48.53 | 31.56      | 43.86 |
| ByT5        | 1.30     | 4.11  | 14.36        | 38.20 | 13.93      | 34.28 |
| ByT5-EU     | 0.64     | 2.30  | 11.36        | 29.74 | 9.80       | 25.04 |
| ByT5-TV     | 1.37     | 3.92  | 3.77         | 11.69 | 3.51       | 9.93  |
| GPT-4-6shot | 2.27     | 5.21  | 9.74         | 21.41 | 8.34       | 25.15 |

mT5 performs worst overall, much lower than ByT5.

# Results (synthetic data)

| Model     | Europarl    |             | Friends (TV) |       | Para5 (TV) |       |
|-----------|-------------|-------------|--------------|-------|------------|-------|
|           | CER         | WER         | CER          | WER   | CER        | WER   |
| AllGreek  | 2.98        | 11.54       | 7.56         | 19.62 | 5.32       | 16.32 |
| RBLSM     | 4.71        | 15.92       | 14.29        | 39.67 | 12.98      | 36.53 |
| RBNLM     | 1.34        | 4.45        | 9.68         | 27.03 | 7.36       | 20.55 |
| mT5       | 6.06        | 10.14       | 32.60        | 48.53 | 31.56      | 43.86 |
| ByT5      | 1.30        | 4.11        | 14.36        | 38.20 | 13.93      | 34.28 |
| ByT5-EU   | <b>0.64</b> | <b>2.30</b> | 11.36        | 29.74 | 9.80       | 25.04 |
| ByT5-TV   | 1.37        | 3.92        | 3.77         | 11.69 | 3.51       | 9.93  |
| T-4-6shot | 2.27        | 5.21        | 9.74         | 21.41 | 8.34       | 25.15 |

mT5 performs worst overall, much lower than ByT5.

ByT5-EU is the best model in its specific domain.

# Results (synthetic data)

| Model     | Europarl    |             | Friends (TV) |              | Para5 (TV)  |             |
|-----------|-------------|-------------|--------------|--------------|-------------|-------------|
|           | CER         | WER         | CER          | WER          | CER         | WER         |
| AllGreek  | 2.98        | 11.54       | 7.56         | 19.62        | 5.32        | 16.32       |
| RBSLM     | 4.71        | 15.92       | 14.29        | 39.67        | 12.98       | 36.53       |
| RBNLM     | 1.34        | 4.45        | 9.68         | 27.03        | 7.36        | 20.55       |
| mT5       | 6.06        | 10.14       | 32.60        | 48.53        | 31.56       | 43.86       |
| ByT5      | 1.30        | 4.11        | 14.36        | 38.20        | 13.93       | 34.28       |
| ByT5-EU   | <b>0.64</b> | <b>2.30</b> | 11.36        | 29.74        | 9.80        | 25.04       |
| ByT5-TV   | 1.37        | 3.92        | <b>3.77</b>  | <b>11.69</b> | <b>3.51</b> | <b>9.93</b> |
| T-4-6shot | 2.27        | 5.21        | 9.74         | 21.41        | 8.34        | 25.15       |

mT5 performs worst overall, much lower than ByT5.

ByT5-EU is the best model in its specific domain.

ByT5-TV is also the best in its domain.

# Results (real-world)

| Model       | Survivorbot |       | Gazzetta (accepted) |       | Gazzetta (rejected) |       |
|-------------|-------------|-------|---------------------|-------|---------------------|-------|
|             | CER         | WER   | CER                 | WER   | CER                 | WER   |
| AllGreek    | 14.66       | 30.48 | 9.71                | 24.54 | 11.17               | 28.94 |
| RBSLM       | 22.70       | 55.16 | 18.55               | 43.01 | 20.70               | 43.07 |
| RBNLM       | 19.99       | 50.21 | 14.59               | 33.10 | 17.41               | 35.42 |
| mT5         | 40.99       | 59.67 | 30.02               | 43.70 | 39.26               | 50.03 |
| ByT5        | 29.45       | 59.89 | 18.23               | 35.01 | 27.63               | 42.49 |
| ByT5-EU     | 22.96       | 51.18 | 15.57               | 32.08 | 24.65               | 37.70 |
| ByT5-TV     | 17.70       | 39.78 | 13.19               | 31.11 | 22.30               | 35.63 |
| ByT5-90shot | 16.41       | 38.22 | 11.17               | 25.98 | 17.69               | 29.61 |
| GPT-4-6shot | 9.44        | 22.74 | 8.02                | 18.36 | 10.80               | 21.76 |



# Results (real-world)

| Model       | Survivorbot |              | Gazzetta (accepted) |              | Gazzetta (rejected) |              |
|-------------|-------------|--------------|---------------------|--------------|---------------------|--------------|
|             | CER         | WER          | CER                 | WER          | CER                 | WER          |
| AllGreek    | 14.66       | 30.48        | 9.71                | 24.54        | 11.17               | 28.94        |
| RBSLM       | 22.70       | 55.16        | 18.55               | 43.01        | 20.70               | 43.07        |
| RBNLM       | 19.99       | 50.21        | 14.59               | 33.10        | 17.41               | 35.42        |
| mT5         | 40.99       | 59.67        | 30.02               | 43.70        | 39.26               | 50.03        |
| ByT5        | 29.45       | 59.89        | 18.23               | 35.01        | 27.63               | 42.49        |
| ByT5-EU     | 22.96       | 51.18        | 15.57               | 32.08        | 24.65               | 37.70        |
| ByT5-TV     | 17.70       | 39.78        | 13.19               | 31.11        | 22.30               | 35.63        |
| ByT5-90shot | 16.41       | 38.22        | 11.17               | 25.98        | 17.69               | 29.61        |
| GPT-4-6shot | <b>9.44</b> | <b>22.74</b> | <b>8.02</b>         | <b>18.36</b> | <b>10.80</b>        | <b>21.76</b> |

GPT-4 is the best performing model in all datasets.

# Results (real-world)

| Model       | Survivorbot |              | Gazzetta (accepted) |              | Gazzetta (rejected) |              |
|-------------|-------------|--------------|---------------------|--------------|---------------------|--------------|
|             | CER         | WER          | CER                 | WER          | CER                 | WER          |
| AllGreek    | 14.66       | 30.48        | 9.71                | 24.54        | 11.17               | 28.94        |
| RBSLM       | 22.70       | 55.16        | 18.55               | 43.01        | 20.70               | 43.07        |
| RBNLM       | 19.99       | 50.21        | 14.59               | 33.10        | 17.41               | 35.42        |
| mT5         | 40.99       | 59.67        | 30.02               | 43.70        | 39.26               | 50.03        |
| ByT5        | 29.45       | 59.89        | 18.23               | 35.01        | 27.63               | 42.49        |
| ByT5-EU     | 22.96       | 51.18        | 15.57               | 32.08        | 24.65               | 37.70        |
| ByT5-TV     | 17.70       | 39.78        | 13.19               | 31.11        | 22.30               | 35.63        |
| ByT5-90shot | 16.41       | 38.22        | 11.17               | 25.98        | 17.69               | 29.61        |
| GPT-4-6shot | <b>9.44</b> | <b>22.74</b> | <b>8.02</b>         | <b>18.36</b> | <b>10.80</b>        | <b>21.76</b> |

GPT-4 is the best performing model in all datasets.

Fine-tuning on limited real world data drastically improves ByT5.

# Results (real-world)

| Model       | Survivorbot |              | Gazzetta (accepted) |              | Gazzetta (rejected) |              |
|-------------|-------------|--------------|---------------------|--------------|---------------------|--------------|
|             | CER         | WER          | CER                 | WER          | CER                 | WER          |
| AllGreek    | 14.66       | 30.48        | 9.71                | 24.54        | 11.17               | 28.94        |
| RBSLM       | 22.70       | 55.16        | 18.55               | 43.01        | 20.70               | 43.07        |
| RBNLM       | 19.99       | 50.21        | 14.59               | 33.10        | 17.41               | 35.42        |
| mT5         | 40.99       | 59.67        | 30.02               | 43.70        | 39.26               | 50.03        |
| ByT5        | 29.45       | 59.89        | 18.23               | 35.01        | 27.63               | 42.49        |
| ByT5-EU     | 22.96       | 51.18        | 15.57               | 32.08        | 24.65               | 37.70        |
| ByT5-TV     | 17.70       | 39.78        | 13.19               | 31.11        | 22.30               | 35.63        |
| ByT5-90shot | 16.41       | 38.22        | 11.17               | 25.98        | 17.69               | 29.61        |
| GPT-4-6shot | <b>9.44</b> | <b>22.74</b> | <b>8.02</b>         | <b>18.36</b> | <b>10.80</b>        | <b>21.76</b> |

AllGreek is the second best model in all real-world experiments.

GPT-4 is the best performing model in all datasets.

Fine-tuning on limited real world data drastically improves ByT5.

# Results (real-world)

Observation: more real-world data (regardless of source) leads to better performance.

| Model       | Survivorbot |              | Gazzetta (accepted) |              | Gazzetta (rejected) |              |
|-------------|-------------|--------------|---------------------|--------------|---------------------|--------------|
|             | CER         | WER          | CER                 | WER          | CER                 | WER          |
| AllGreek    | 14.66       | 30.48        | 9.71                | 24.54        | 11.17               | 28.94        |
| RBSLM       | 22.70       | 55.16        | 18.55               | 43.01        | 20.70               | 43.07        |
| RBNLM       | 19.99       | 50.21        | 14.59               | 33.10        | 17.41               | 35.42        |
| mT5         | 40.99       | 59.67        | 30.02               | 43.70        | 39.26               | 50.03        |
| ByT5        | 29.45       | 59.89        | 18.23               | 35.01        | 27.63               | 42.49        |
| ByT5-EU     | 22.96       | 51.18        | 15.57               | 32.08        | 24.65               | 37.70        |
| ByT5-TV     | 17.70       | 39.78        | 13.19               | 31.11        | 22.30               | 35.63        |
| ByT5-90shot | 16.41       | 38.22        | 11.17               | 25.98        | 17.69               | 29.61        |
| GPT-4-6shot | <b>9.44</b> | <b>22.74</b> | <b>8.02</b>         | <b>18.36</b> | <b>10.80</b>        | <b>21.76</b> |

AllGreek is the second best model in all real-world experiments.

GPT-4 is the best performing model in all datasets.

Fine-tuning on limited real world data drastically improves ByT5.

# Error Analysis

|                                  |           |  |
|----------------------------------|-----------|--|
| <i>other-language characters</i> | GREEKLISH | gia na to doume to allani to keno mprogiobits an kalucei               |
|                                  | GOLD      | Για να το δούμε το αλάιι το κενό Βοιονιc αν καλύψει                    |
|                                  | MT5       | για να το δούμε το άλνενο το κένο μπορντι <del>ρ</del> αईΣ αν καλύπτει |
| <i>incorrect spelling</i>        | GREEKLISH | to gipedo tis aek to eidane? :P  |
|                                  | GOLD      | Το γήπεδο της ΑΕΚ το είδανε; :P  |
|                                  | BYT5      | Το γυπέδο της αεκ το είδανε? :P  |
|                                  | MT5       | Το γύρω πόλων της άμες το έθατε?                                       |
| <i>non-transliterated words</i>  | GREEKLISH | ginei ke files, boris na mu peis;                                      |
|                                  | GOLD      | γίνει και φίλες, μπορείς να μου πεις;                                  |
|                                  | ALLGREEK  | γίνει και φίλες, boris να μου πεις;                                    |

# Error Analysis

|                                  |           |  |
|----------------------------------|-----------|--|
| <i>other-language characters</i> | GREEKLISH | gia na to doume to allani to keno mprogiobits an kalucei                           |
|                                  | GOLD      | Για να το δούμε το αλάني το κενό Βοιονιc αν καλύψει                                |
|                                  | MT5       | για να το δούμε το άλνενο το κένο μπροντ <del>ι</del> <sup>रअई</sup> ς αν καλύπτει |
| <i>incorrect spelling</i>        | GREEKLISH | to gipedo tis aek to eidane? :P  |
|                                  | GOLD      | Το γήπεδο της ΑΕΚ το είδανε; :P  |
|                                  | BYT5      | Το γυπέδο της αεκ το είδανε? :P  |
|                                  | MT5       | Το γύρω πόλων της άμες το έθατε?   |
| <i>non-transliterated words</i>  | GREEKLISH | ginei ke files, boris na mu peis,  |
|                                  | GOLD      | γίνει και φίλες, μπορείς να μου πεις   |
|                                  | ALLGREEK  | γίνει και φίλες, boris να μου πεις;  |

mT5's multilingual sub-word tokens might not be optimal for Greeklish.

# Error Analysis

|                                  |           |   |
|----------------------------------|-----------|---|
| <i>other-language characters</i> | GREEKLISH | gia na to doume to allani to keno mprogiobits an kalucei          |
|                                  | GOLD      | Για να το δούμε το αλάι το κενό Βοιονιc αν καλύπει                |
|                                  | MT5       | για να το δούμε το άλνενο το κένο μπορντ <sup>Α</sup> αν καλύπτει |
| <i>incorrect spelling</i>        | GREEKLISH | to gipedo tis aek to eidane? :P                                   |
|                                  | GOLD      | Το γήπεδο της ΑΕΚ το είδανε; :P                                   |
|                                  | BYT5      | Το γυπέδο της αεκ το είδανε? :P                                   |
|                                  | MT5       | Το γύρω πόλων της άμες το έθατε?                                  |
| <i>non-transliterated words</i>  | GREEKLISH | ginei ke files, boris na mu peis,                                 |
|                                  | GOLD      | γίνει και φίλες, μπορείς να μου πεις                              |
|                                  | ALLGREEK  | γίνει και φίλες, boris να μου πεις;                               |

When AllGreek comes across input that it struggles with, it often skips transliterating it (maybe identifying it as code-switching).

mT5's multilingual sub-word tokens might not be optimal for Greeklish.

# Terms altered by GPT-based

| Greeklsh (SOURCE)  | Greek (TRANSLITERATED)  |
|--|---|
| An <b>i united</b> eprepe na valei <b>k</b> 4o gia na perasei tha to evaze xalara  | Αν <b>η United</b> έπρεπε να βάλει <b>και</b> 4ο για να περάσει, θα το έβαζε άνετα  |
| Ta topika einai ola xera kai <b>giafto</b> den paei kanena paidaki na spasi ta podia tou   | Τα τοπικά είναι όλα ξερά και <b>γι'αυτό</b> δεν πάει κανένα παιδάκι να σπάσει τα πόδια του  |
| Na afairethoun <b>t</b> vraveia apo ton <b>berg</b> ton <b>anastasiou</b> kai ton <b>risvani</b> na <b>t</b> dwsoun ston <b>mitroglou p</b> einai stin agglia ston <b>mitsel</b> pou pire prwtathlima apo ton <b>au-gousto p</b> perase kai kalutero neo paixth na to dwsoun ston <b>vergo</b> | Να αφαιρεθούν <b>τα</b> βραβεία από τον <b>Βεργκ</b> τον <b>Αναστασίου</b> και τον <b>Ρισβάνη</b> να <b>τα</b> δώσουν στον <b>Μήτρογλου που</b> είναι στην <b>Αγγλία</b> στον <b>Μίτσελ</b> που πήρε πρωτάθλημα από τον <b>Αύγουστο που</b> πέρασε και για καλύτερο νέο παίχτη να το δώσουν στον <b>Βέργο</b> |

- GPT-based models try to clarify **abbreviations**, **short forms**, and **named entities**.
- Usually leads to less noisy & better transliteration.



# Downstream Task: Content Moderation

Instruction  
prompt

*You are an expert in online content moderation. The following is an online user comment in Greek or Greeklisch, which can either be accepted for publication or rejected as toxic by the moderator. Classify the next candidate post as either 1 for “rejected” or 0 for “accepted”*



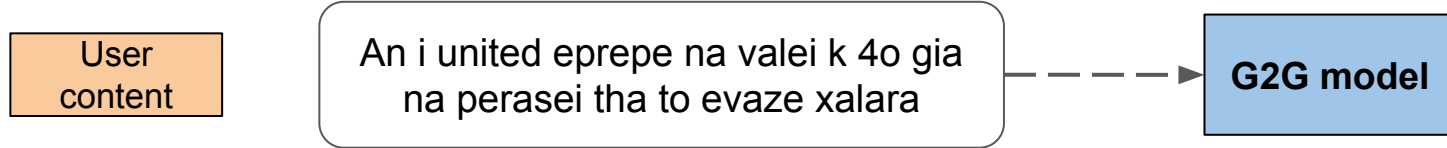
User  
content

An i united eprepe na valei k 4o gia  
na perasei tha to evaze xalara

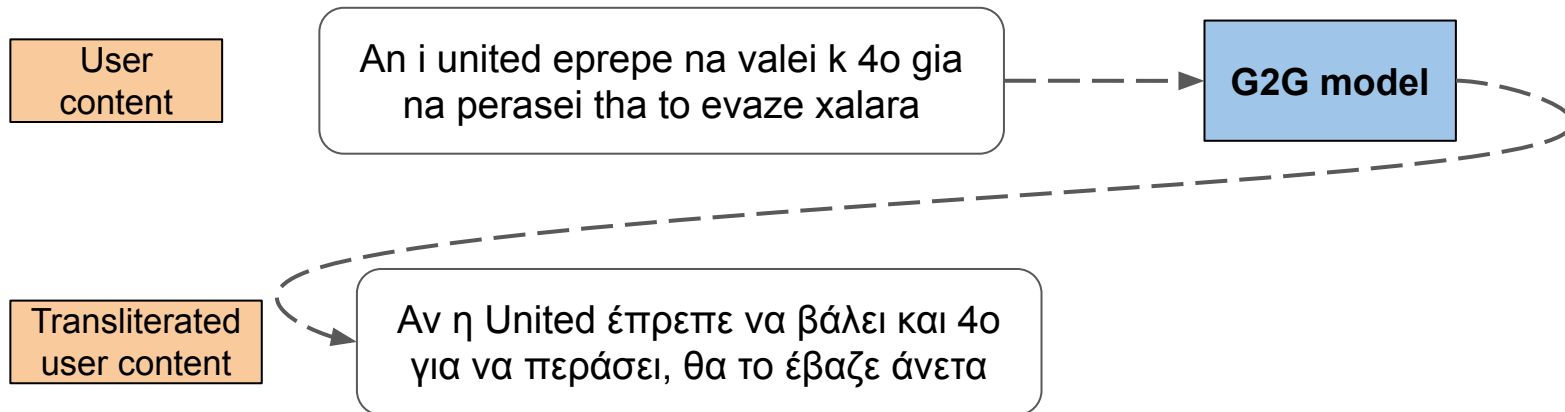
GPT

Accept / Reject

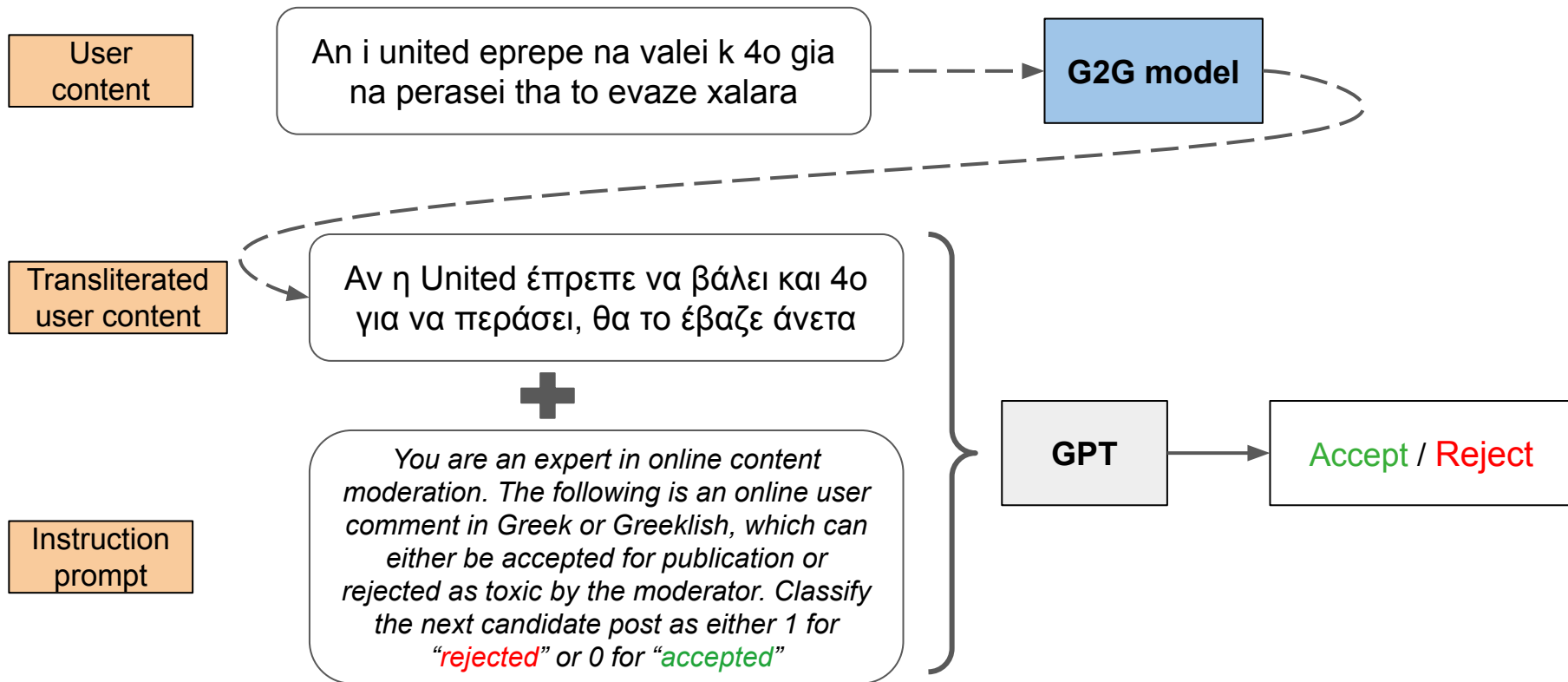
# Downstream Task: Content Moderation (with transliteration)



# Downstream Task: Content Moderation (with transliteration)



# Downstream Task: Content Moderation (with transliteration)



# Downstream Task: Content Moderation

|   | Greeklish        |               | Transliterated   |               |
|---|------------------|---------------|------------------|---------------|
| <b>Accuracy</b>   | 0.49             |               | <b>0.53</b>      |               |
| <b>F<sub>2</sub> (P<sub>rejected</sub>, P<sub>accepted</sub>)</b> | 0.54             |               | <b>0.69</b>      |               |
|   | <b>Precision</b> | <b>Recall</b> | <b>Precision</b> | <b>Recall</b> |
| <b>Accepted</b>   | 0.47             | 0.71          | <b>0.50</b>      | <b>0.89</b>   |
| <b>Rejected</b>   | 0.54             | <b>0.30</b>   | <b>0.69</b>      | 0.21          |
| <b>Average</b>  | 0.51             | 0.51          | <b>0.59</b>      | <b>0.55</b>   |

- On average, adding a transliteration step helps (even if it's not perfect).
- No change needed in the model of the downstream task!
- We use F<sub>2</sub> as defined by Pavlopoulos et al. (2017) (<https://aclanthology.org/D17-1117.pdf>)

# Conclusions

- We evaluated various approaches for G2G transliteration, featuring approaches based on:
  - Transliteration Rules + {Statistical / Deep Learning} LM
  - Fine-tuned multilingual transformers
  - A previous system based on statistical machine learning
  - Zero/Zero/Few-shot prompting LLMs (GPT-3.4, GPT-4)

# Conclusions

- We evaluated various approaches for G2G transliteration, featuring approaches based on:
  - Transliteration Rules + {Statistical / Deep Learning} LM
  - Fine-tuned multilingual transformers
  - A previous system based on statistical machine learning
  - Zero/Zero/Few-shot prompting LLMs (GPT-3.4, GPT-4)
- **We conclude that:**

# Conclusions

- We evaluated various approaches for G2G transliteration, featuring approaches based on:
  - Transliteration Rules + {Statistical / Deep Learning} LM
  - Fine-tuned multilingual transformers
  - A previous system based on statistical machine learning
  - Zero/Few-shot prompting LLMs (GPT-3.4, GPT-4)
- **We conclude that:**
- Zero/Few-shot prompting LLMs yields the best results.  
..but..
- An old system based of statistical LM can be competitive, often outperforming newer approaches and techniques.
- Greeklish-to-Greek transliteration is far from solved.  
..but..
- Even non-perfect transliteration can help significantly in downstream tasks (experimental evidence from content moderation).



# Thank you!



<https://github.com/nlpauieb/greeklish>



@AUEBNLPGroup



{touzatos,ion,annis}@aueb.gr