# Patterns of misspellings in L2 English – a view from the ETS Spelling Corpus

Michael Flor,
Yoko Futagi,
Melissa Lopez,
Matthew Mulholland

NLP & Speech Group, R&D Division
Educational Testing Service, Princeton, NJ, USA

mflor@ets.org

# The Corpus

- **Initial Motivation:**
  evaluate speller performance
  by comparing it to human-annotated data.

- **We needed:**
  a human-annotated corpus of misspellings,
  where misspellings appear in their original context.

- **Method:**
  Use ConSpel to generate a corpus automatically,
  then let human annotators work on it…

# The Corpus

**Materials** – English essays written on TOEFL and GRE tests at international testing centers around the world. (computer-based delivery, QWERTY keyboard),

| Program/task | Description of writing activity |
|---|---|
| TOEFL Independent | support an *opinion* in writing (topic assigned). |
| TOEFL Integrated | write essay responses based on reading and listening tasks (**summarize and compare arguments**) |
| GRE Issue | express *opinion* clearly, in writing, about a topic of general interest (topic assigned). |
| GRE Argument | analyze and evaluate **arguments** according to specific instructions and convey evaluation clearly in writing. |

# The Corpus

- 4 program/task groups
- 10 different prompts for each task
- 75 essays per prompt
- Total:   3,000 essays (963K words)
- Essay length ranges from 29 to798 words, average 321 words

# Annotation software

# Annotation software

- Each essay was annotated with automated support
- Each misspeling was marked and provided with a <u>correction</u>
- Software automatically defined type of misspelling
- Full essay context is seen during annotation
- Essay and annotation are stored in XML with stand-off markup

# The Corpus

Inter-Annotator Agreement

- Each essay was annotated by two annotators.
- Annotators strictly agreed in 82.6% the cases.
- Inter-annotator agreement was calculated over all words of the corpus: 99.3%.
- Cohen's Kappa=0.85, p<0.001.
- All differences and difficulties were resolved by a third annotator (adjudicator).

# Types and counts of misspellings

| | Description | Count in corpus |
|---|---|---|
| 1 | single token non-word misspelling (e.g. "**businees**") also includes fusion errors (e.g. "**niceday**" for "*nice day*") | 21142 (**80.05%**) |
| 2 | misspelling (?) (non-word token for which no plausible correction was found) | 52 (**0.20%**) |
| 3 | single token real-word misspelling (e.g. "**they**" for "*then*") | 3393 (**12.85%**) |
| 4 | multi-token misspelling with at kleast one non-word (e.g. "**mor efun**" for "*more fun*") | 574 (**2.17%**) |
| 5 | multi-token real-word misspelling (e.g. "**with out**" for "*without*") | 1251 (**4.73%**) |
| | Total | 26412 (**100%**) |

# Breakdown by program/task

| | GRE Argument | GRE Issue | TOEFL Independent | TOEFL Integrated | TOTAL |
|---|---|---|---|---|---|
| Essays | 750 | 750 | 750 | 750 | 3,000 |
| Without misspellings | 60 | 21 | 18 | 21 | 120 |
| Word Count | 263,578 | 336,301 | 212,930 | 151,031 | 963,840 |
| Average WC | *351* | *448* | *284* | *201* | *321* |
| **Misspellings** | 5935 | 7962 | 7285 | 5230 | 26412 |
| % of all words | 2.25% | 2.37% | 3.42% | 3.46% | 2.74% |

# Breakdown by error-type and program/task

# Count of NS/NNS essays

| | TOEFL | GRE | Total count | Essays without misspellings |
|---|---|---|---|---|
| NS | 19 | 634 | 653 | 67 **(10.7%)** |
| NNS | 1481 | 866 | 2347 | 53 **(2.3%)** |

Corpus GRE essays by score and NS-status

Non-native speakers of English (ELLs) are more prone to making spelling errors ?

Consider proficiency

# Spelling Error density



Average % of misspelled words per essay, by NS/NNS and score, GRE only (1500 essays)

- For each population, average percent of misspelled words (per essay) <u>decreases</u> with better proficiency

- There is a gap between NS & NNS at lower proficiencies, (native English speakers make less misspellings, on average) but the gap is closing 'quickly' ! (both main effects and interaction are sig., p<.0001)

# How often is the first character different?

| | total | 1st diff | % |
|---|---|---|---|
| Non-word | 21142 | 522 | 2.47% |
| Real word | 3393 | 404 | 11.91% |
| Multi-token | 574 | 10 | 1.74% |
| Multi-token RW | 1251 | 7 | 0.56% |

Breakdown by NS/NNS

| | | total | 1st diff | % |
|---|---|---|---|---|
| Non-word | NNS | 18264 | 465 | 2.50% |
| | NS | 2878 | 57 | 1.98% |
| Real word | NNS | 3008 | 361 | 12.00% |
| | NS | 385 | 43 | 11.17% |

## Examples

Non-words

> **onformation** information
> **imerged** emerged
> **onther** another
> **htis** this
> **phorensic** forensic
> **tasttime** pasttime

Real words

> **write** right
> **equality** quality
> **asocial** social
> **affect** effect
> **participated** anticipated
> **as** has

| Dist. (LED) | Total tokens | Count NS | % NS | Count NNS | % NNS |
|---|---|---|---|---|---|
| 1 | 16908 | 2393 | 83.15% | 14515 | 79.47% |
| 2 | 2957 | 372 | 12.93% | 2585 | 14.15% |
| 3 | 827 | 88 | 3.06% | 739 | 4.05% |
| 4 | 296 | 22 | 0.76% | 274 | 1.50% |
| 5 | 100 | 2 | 0.07% | 98 | 0.54% |
| 6 | 41 | 1 | 0.03% | 40 | 0.22% |
| 7 | 7 | | | 7 | 0.04% |
| 8 | 2 | | | 2 | 0.01% |
| 9 | 4 | | | 4 | 0.02% |

```
recom recommendation (9)
unsatisfy dissatisfaction (9)
naiberhouad neighborhood (6)
chraterics characteristics (5)
voultaneer volunteer (4)
metirals materials (3)
```

The difference
83.1% vs. 79.4%
is significant (p <.0001),
but misleading

GRE data:
significant main effect of Score ($p<0.001$),
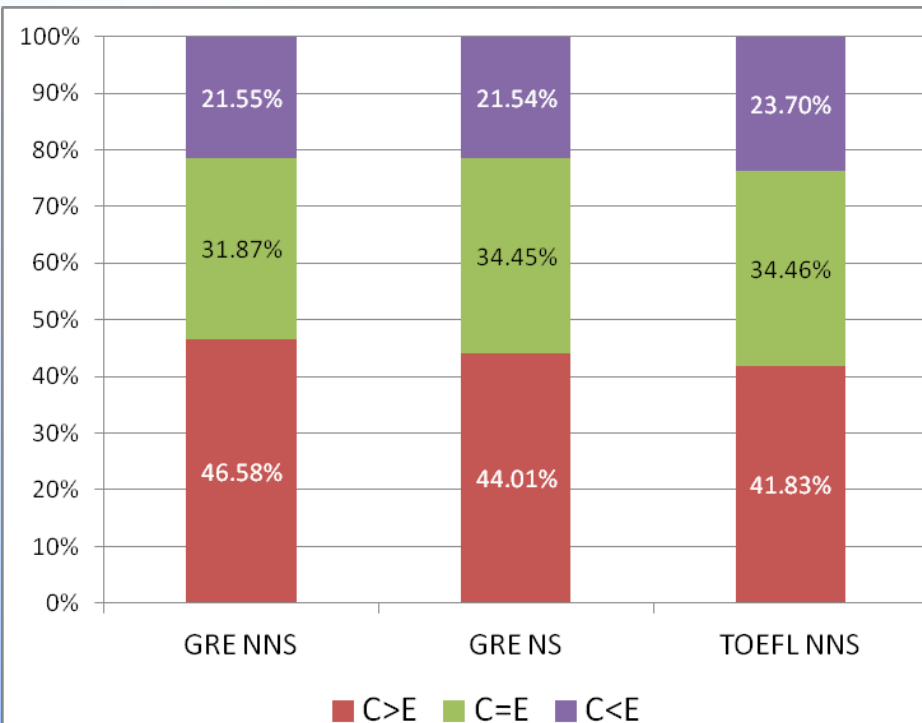no effect of NS/NNS, ($p=0.38$)
and no interaction ($p=0.155$).

TOEFL data:
significant  effect of Score ($p<0.001$).

For 1-token NW errors,
'severity of error' (DLED) depends on proficiency, not NS/NNS distinction;
and yet…

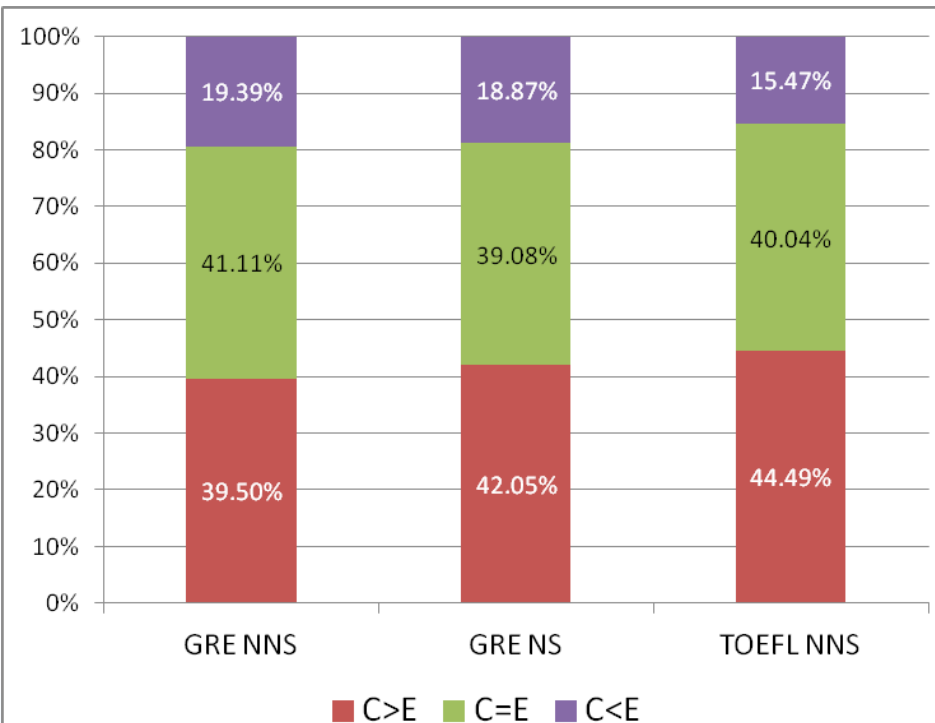# Length of error-form vs. correct-form



**1-token NW  n=21059**

8004          2795          10260
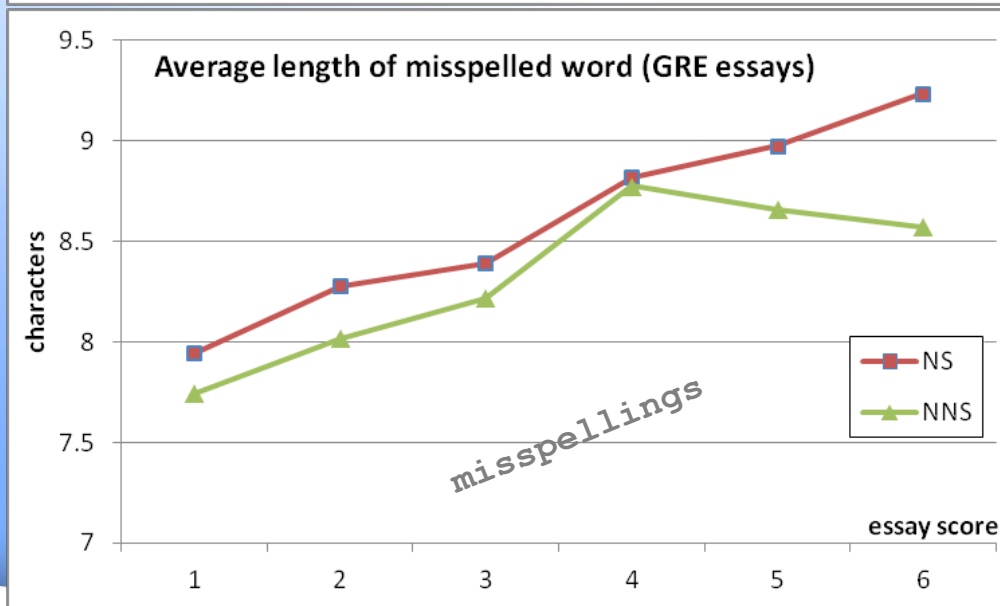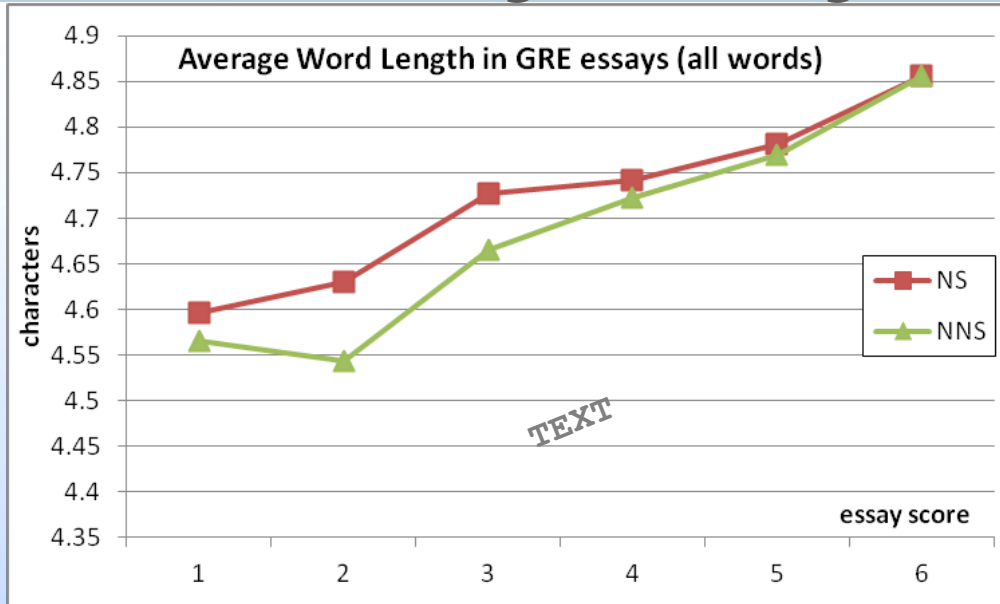
**1-token RW  n=3379**

1547          371          1461

For 1-token NW errors,  and for 1-token RW errors:

For all groups, when a word is misspelled,
there is a tendency to 'miss' characters, rather than to 'add' characters!
And a strong tendency to preserve length!

```
Onformation (=) information
as (<) has
asocial (>) social
```
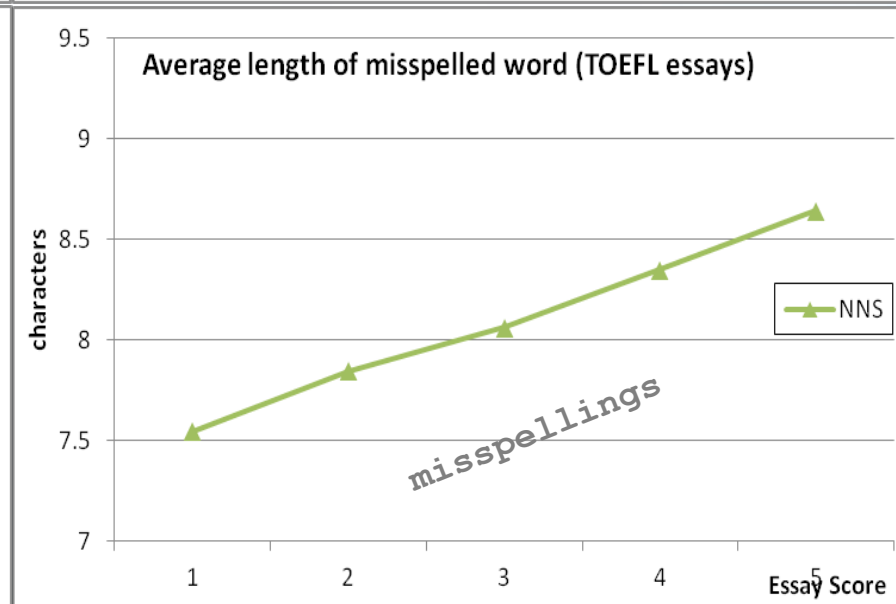
# Average word length and spelling (1-token NW)
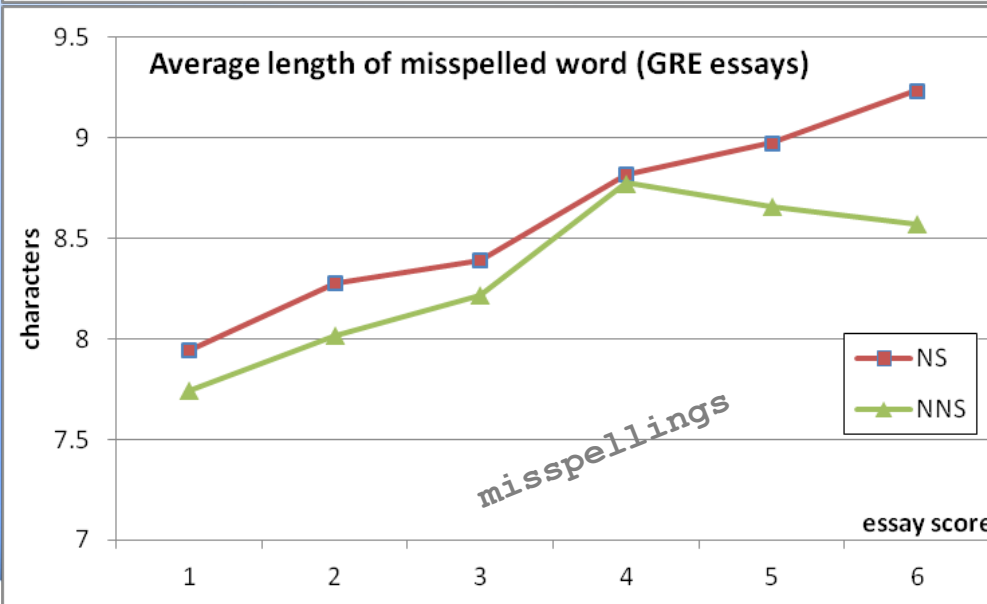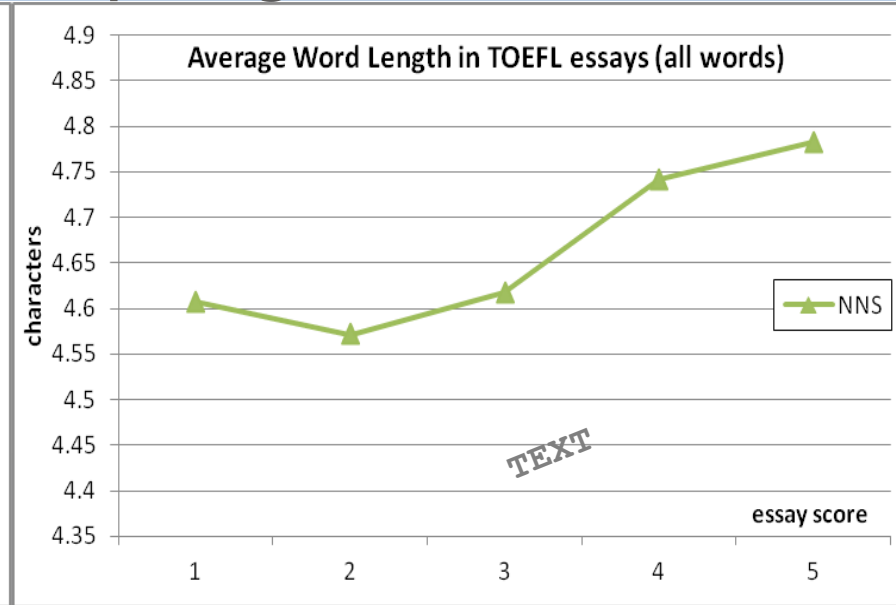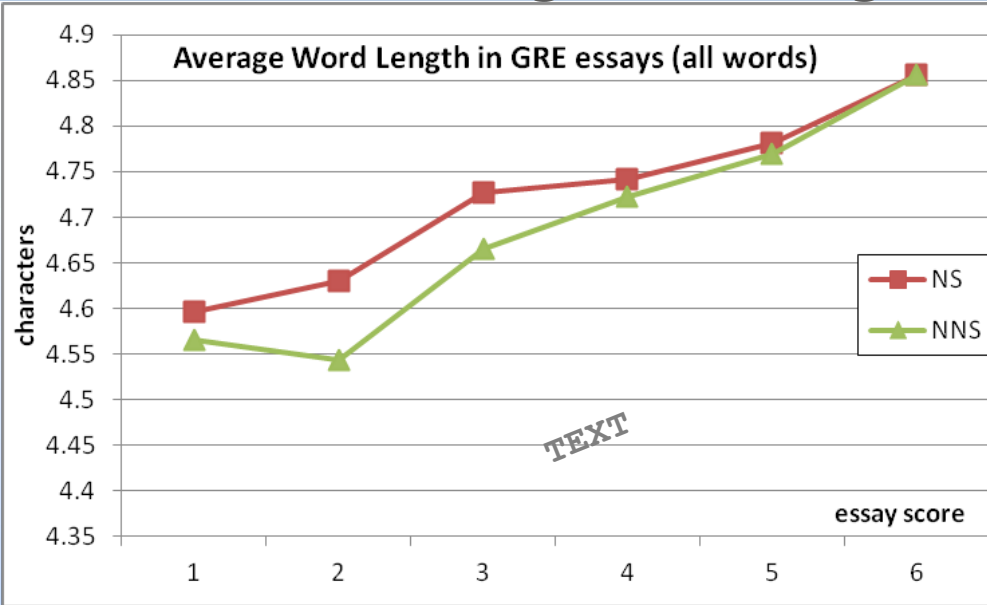


**1500 GRE essays**

- Average word length (per essay) <u>increases</u> with better proficiency.

- NS typically use more long words

- The gap is rapidly closing with better proficiency
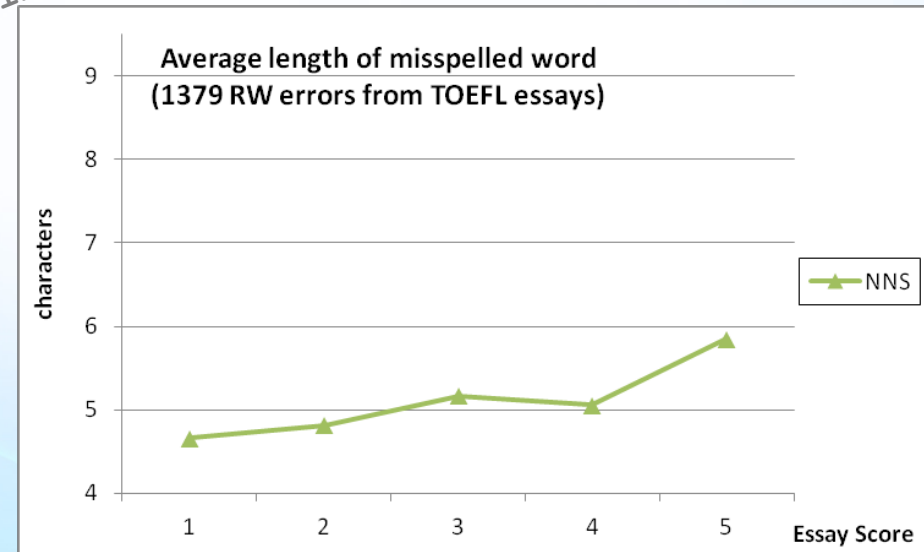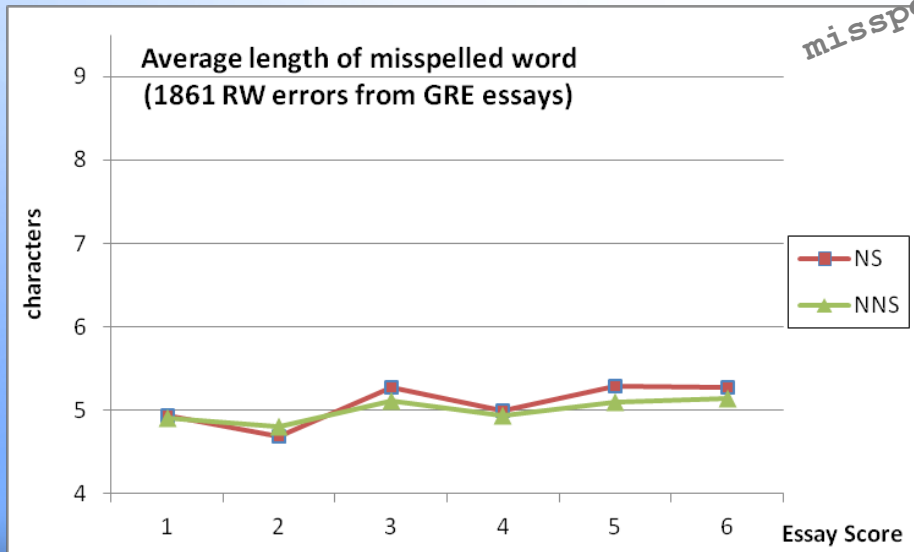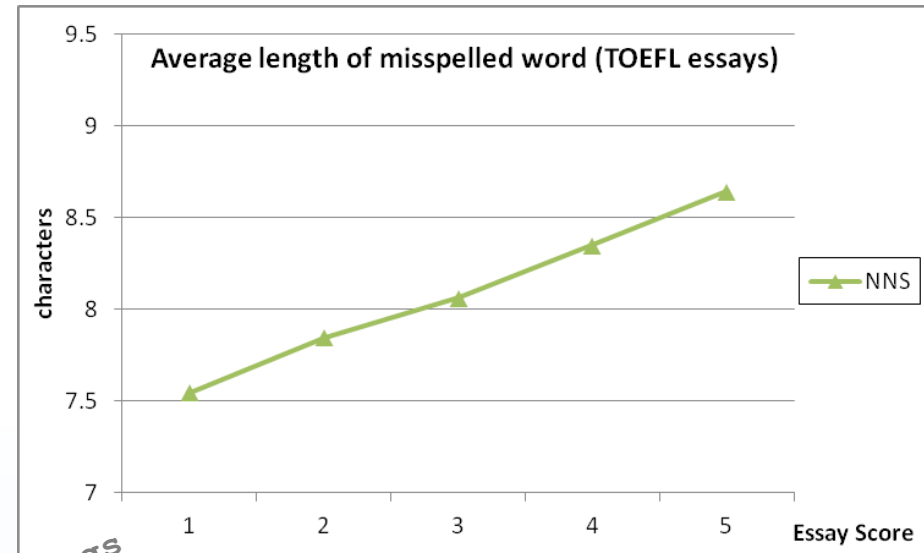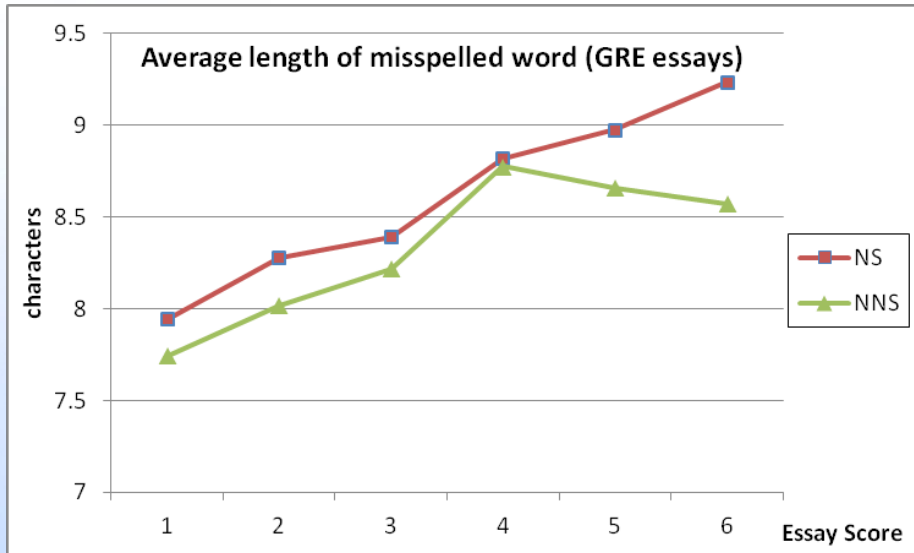
- (both main effects and interaction are sig., p<.0001)

**10110 1-token NW errors (GRE essays)**

- Average length of intended word (misspelled to NW) <u>increases</u> with better proficiency.

- NS typically err in the longer words

- The gap closes at score=4, then widens!
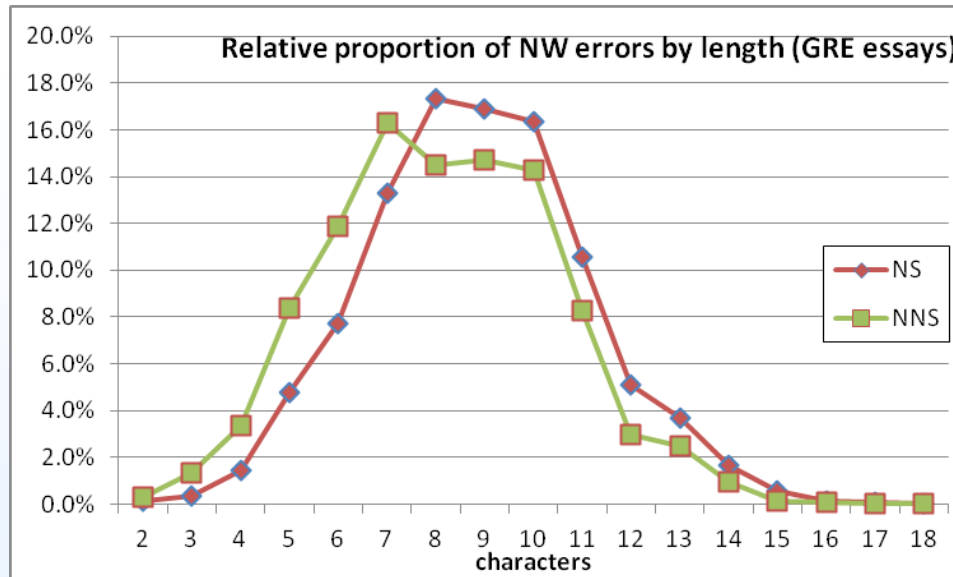
- (both main effects and interaction are sig., p<.0001)
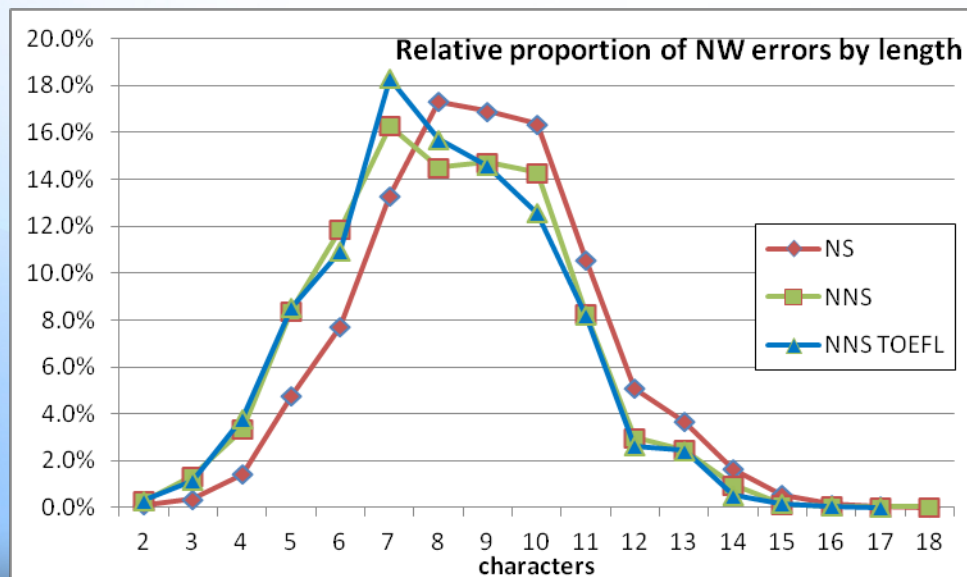
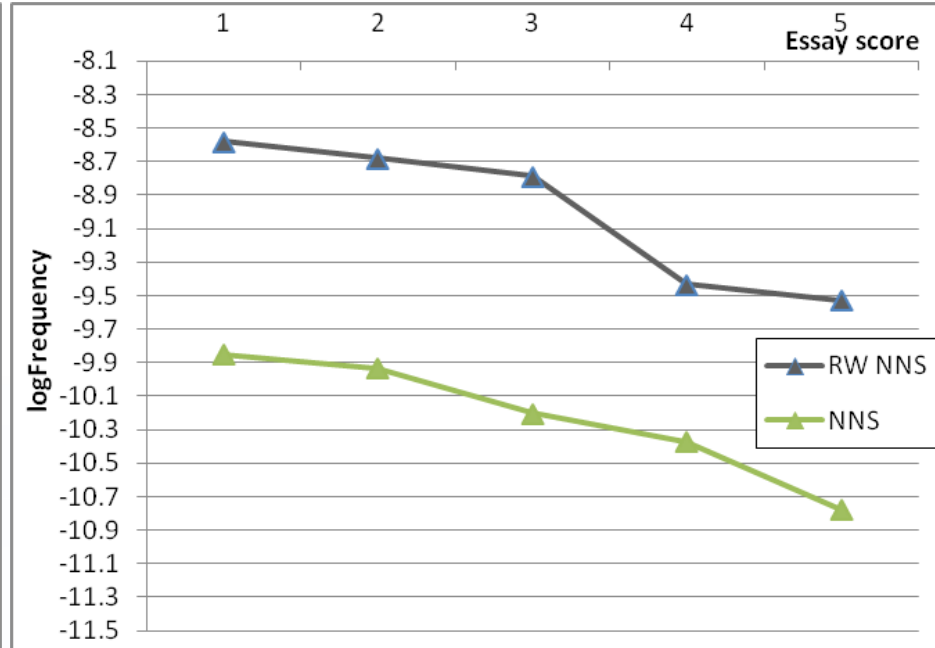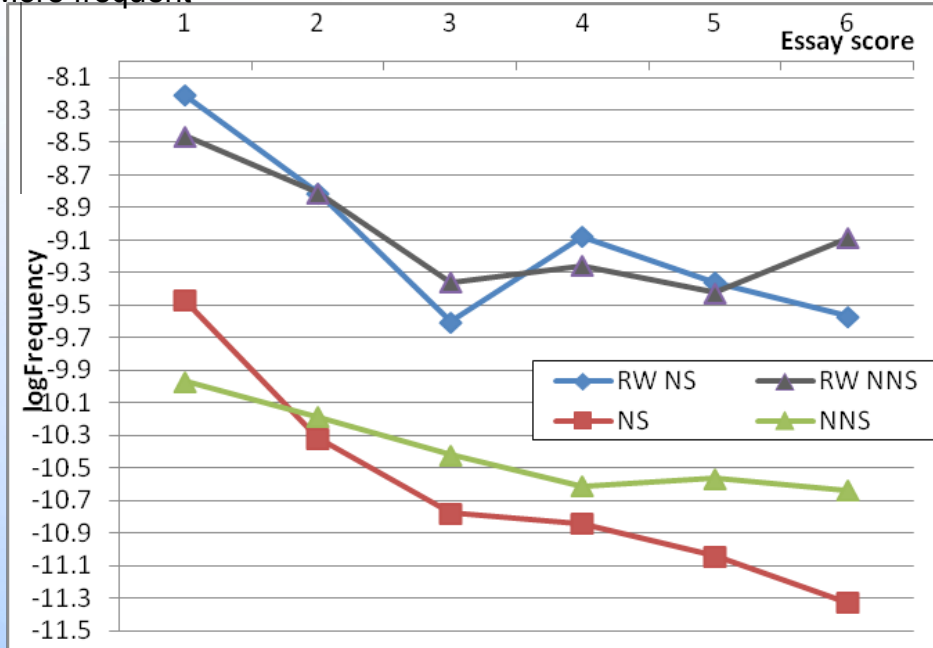# Average word length and spelling (1-token NW)

1500 GRE essays

+1500 TOEFL essays

# Word frequency and spelling (1-token NW & RW)

GRE data

TOEFL data

More frequent



logFrequency of the corrected-form of a misspelling

onformation information

For 1-token NW errors, GRE data: both main effects and interaction are sig., p<.002.
For 1-token RW errors, GRE data: no effect is sig. (even Score p=0.71).
TOEFL data, for each NW and RW: effect of Score is sig., p<.001.

The differences between NW and RW are sig. (p<.001) in each of 3 comparisons:
The average frequency of words where RW errors are made is higher than
average frequency of words where NW errors are made.

# Thank You

mflor@ets.org