# Lecture 2: Writers' aids: Spelling errors

LING-351 Language Technology and LLMs

Instructor: Hakyung Sung

August 28, 2025

# Table of contents

1

# Review

- Language

# Language, writing, encoding

- Language
- Writing

- Language
- Writing
- Language = writing?

# Language, writing, encoding

- Language
- Writing
- Language = writing?
- Three major systems to encode languages:

- Language
- Writing
- Language = writing?
- Three major systems to encode languages:
  - Alphabetic

- Language
- Writing
- Language = writing?
- Three major systems to encode languages:
  - Alphabetic → sounds

- Language
- Writing
- Language = writing?
- Three major systems to encode languages:
  - Alphabetic $\rightarrow$ sounds
  - Syllabic

- Language
- Writing
- Language = writing?
- Three major systems to encode languages:
  - Alphabetic → sounds
  - Syllabic → syllables

- Language
- Writing
- Language = writing?
- Three major systems to encode languages:
    - Alphabetic $\rightarrow$ sounds
    - Syllabic $\rightarrow$ syllables
    - Logographic

- Language
- Writing
- Language = writing?
- Three major systems to encode languages:
    - Alphabetic $\rightarrow$ sounds
    - Syllabic $\rightarrow$ syllables
    - Logographic $\rightarrow$ meanings

- Digital technology can be understood as another form of writing to encode language into digital formats

- Digital technology can be understood as another form of writing to encode language into digital formats
- Bit (0/1 signal): the smallest unit of digital information

- Digital technology can be understood as another form of writing to encode language into digital formats
- Bit (0/1 signal): the smallest unit of digital information
- Byte (8 bits): a bundle of 8 bits, the basic unit of storage

- Digital technology can be understood as another form of writing to encode language into digital formats
- Bit (0/1 signal): the smallest unit of digital information
- Byte (8 bits): a bundle of 8 bits, the basic unit of storage
- Character encoding (UTF-8): rules that map bytes to code points

# Lesson plan

- ~~Review~~
- Spelling problems in writing

- ~~Review~~
- Spelling problems in writing
- Different types of spelling errors

- ~~Review~~
- Spelling problems in writing
- Different types of spelling errors
- Building a simple spell-checker

- ~~Review~~
- Spelling problems in writing
- Different types of spelling errors
- Building a simple spell-checker
- Thinking about a more complex spell-checker

- ~~Review~~
- Spelling problems in writing
- Different types of spelling errors
- Building a simple spell-checker
- Thinking about a more complex spell-checker
- Wrap-up

Key idea: ~~Spelling errors are annoying~~

- ~~Review~~
- Spelling problems in writing
- Different types of spelling errors
- Building a simple spell-checker
- Thinking about a more complex spell-checker
- Wrap-up

Key idea: ~~Spelling errors are annoying~~
Spelling errors vary by types (and even by languages);
there is no one-size-fits-all solution.

# Spelling problems in writing

A-M-E-L-I-O-R-A-T-E.

- English has used the Latin alphabet since the 9th century

# Writing English: A long evolution

- English has used the Latin alphabet since the 9th century
  - Before that, Old English was written in the runic script

## Writing English: A long evolution

- English has used the Latin alphabet since the 9th century
  - Before that, Old English was written in the runic script
  - Christian missionaries introduced the Latin alphabet in the 7th century

- English has used the Latin alphabet since the 9th century
  - Before that, Old English was written in the runic script
  - Christian missionaries introduced the Latin alphabet in the 7th century
  - By the 9th century, it became the dominant writing system

## Writing English: A long evolution

- English has used the Latin alphabet since the 9th century
  - Before that, Old English was written in the runic script
  - Christian missionaries introduced the Latin alphabet in the 7th century
  - By the 9th century, it became the dominant writing system
- Related technologies:

- English has used the Latin alphabet since the 9th century
  - Before that, Old English was written in the runic script
  - Christian missionaries introduced the Latin alphabet in the 7th century
  - By the 9th century, it became the dominant writing system
- Related technologies:
  - Handwriting on parchment and paper (monastic scribes, medieval manuscripts)

- English has used the Latin alphabet since the 9th century
  - Before that, Old English was written in the runic script
  - Christian missionaries introduced the Latin alphabet in the 7th century
  - By the 9th century, it became the dominant writing system

- Related technologies:
  - Handwriting on parchment and paper (monastic scribes, medieval manuscripts)
  - Printing press (1470s in Englan)
    → wider literacy, circulation of books

- English has used the Latin alphabet since the 9th century
  - Before that, Old English was written in the runic script
  - Christian missionaries introduced the Latin alphabet in the 7th century
  - By the 9th century, it became the dominant writing system
- Related technologies:
  - Handwriting on parchment and paper (monastic scribes, medieval manuscripts)
  - Printing press (1470s in Englan)
    → wider literacy, circulation of books
  - Typewriter (1860s)
    → faster, more uniform writing

# Writing English: A long evolution

- English has used the Latin alphabet since the 9th century
  - Before that, Old English was written in the runic script
  - Christian missionaries introduced the Latin alphabet in the 7th century
  - By the 9th century, it became the dominant writing system

- Related technologies:
  - Handwriting on parchment and paper (monastic scribes, medieval manuscripts)
  - Printing press (1470s in Englan)
    → wider literacy, circulation of books
  - Typewriter (1860s)
    → faster, more uniform writing
  - Digital word processing (20th century)
    → autocorrect, spell checkers

# Writing English: A long evolution

- English has used the Latin alphabet since the 9th century
  - Before that, Old English was written in the runic script
  - Christian missionaries introduced the Latin alphabet in the 7th century
  - By the 9th century, it became the dominant writing system
- Related technologies:
  - Handwriting on parchment and paper (monastic scribes, medieval manuscripts)
  - Printing press (1470s in Englan)
    → wider literacy, circulation of books
  - Typewriter (1860s)
    → faster, more uniform writing
  - Digital word processing (20th century)
    → autocorrect, spell checkers
- Standardized spelling came much later...

- Spelling wasn't standardized until the mid-1600s to 1700s

- Spelling wasn't standardized until the mid-1600s to 1700s
- Influenced by:

- Spelling wasn't standardized until the mid-1600s to 1700s
- Influenced by:
  - **King James Bible (1611)** – named after King James I of England, who authorized a new translation by scholars to unify religious practices

- Spelling wasn't standardized until the mid-1600s to 1700s
- Influenced by:
  - **King James Bible (1611)** – named after King James I of England, who authorized a new translation by scholars to unify religious practices
  - **Early dictionaries** – e.g., Robert Cawdrey's *Table Alphabeticall* (1604), Samuel Johnson's dictionary (1755)

- Spelling wasn't standardized until the mid-1600s to 1700s
- Influenced by:
  - **King James Bible (1611)** – named after King James I of England, who authorized a new translation by scholars to unify religious practices
  - **Early dictionaries** – e.g., Robert Cawdrey's *Table Alphabeticall* (1604), Samuel Johnson's dictionary (1755)
- Authors themselves didn't use consistent spelling

- Spelling wasn't standardized until the mid-1600s to 1700s
- Influenced by:
  - **King James Bible (1611)** – named after King James I of England, who authorized a new translation by scholars to unify religious practices
  - **Early dictionaries** – e.g., Robert Cawdrey's *Table Alphabeticall* (1604), Samuel Johnson's dictionary (1755)
- Authors themselves didn't use consistent spelling
- Shakespeare's name appeared in many forms:

## English without standard spelling

- Spelling wasn't standardized until the mid-1600s to 1700s
- Influenced by:
    - **King James Bible (1611)** – named after King James I of England, who authorized a new translation by scholars to unify religious practices
    - **Early dictionaries** – e.g., Robert Cawdrey's *Table Alphabeticall* (1604), Samuel Johnson's dictionary (1755)
- Authors themselves didn't use consistent spelling
- Shakespeare's name appeared in many forms:

- Spelling wasn't standardized until the mid-1600s to 1700s
- Influenced by:
  - **King James Bible (1611)** – named after King James I of England, who authorized a new translation by scholars to unify religious practices
  - **Early dictionaries** – e.g., Robert Cawdrey's *Table Alphabeticall* (1604), Samuel Johnson's dictionary (1755)
- Authors themselves didn't use consistent spelling
- Shakespeare's name appeared in many forms:

*Willm Shakp, William Shaksper, Wm Shakspe, William Shakspere, Willm Shakspere, William Shakspeare*

- Even without standard spelling, we understand:

- Even without standard spelling, we understand:

  *To what extent do the **spellling errers** in this **setnence dirsupt** your **undertsanding**?*

# Why *standardized* spelling?

- Even without standard spelling, we understand:
  *To what extent do the **spelling errers** in this **setnence dirsupt** your **undertsanding**?*

- Readers often focus on **word shape**, not letter-by-letter decoding

- In Shakespeare's time, spelling was flexible.

- In Shakespeare's time, spelling was flexible.
- Imagine replacing English spelling with IPA (phonetic spelling).

- In Shakespeare's time, spelling was flexible.
- Imagine replacing English spelling with **IPA** (phonetic spelling).

**Question**

What are the benefits and drawbacks of having a standardized spelling system?

- Supports literacy across dialects/various pronunciations (e.g., *tomato*, *Atlanta*)

- Supports literacy across dialects/various pronunciations (e.g., *tomato*, *Atlanta*)
- Enables searching and record-keeping

- English spelling remains complex

- English spelling remains complex
- Thankfully, we have **writers' aids**:

- English spelling remains complex
- Thankfully, we have **writers' aids**:
    - Spell checkers

# Spelling in the digital age

- English spelling remains complex
- Thankfully, we have **writers' aids**:
    - Spell checkers
    - Predictive text; auto-complete

- English spelling remains complex
- Thankfully, we have **writers' aids**:
  - Spell checkers
  - Predictive text; auto-complete
  - Generative LLMs

Spelling in the digital age

- English spelling remains complex
- Thankfully, we have **writers' aids**:
    - Spell checkers
    - Predictive text; auto-complete
    - Generative LLMs

Group discussion

- English spelling remains complex
- Thankfully, we have **writers' aids**:
    - Spell checkers
    - Predictive text; auto-complete
    - Generative LLMs

**Group discussion**

- *(Put into the shared deck)* Come up with at least one example

- English spelling remains complex
- Thankfully, we have **writers' aids**:
    - Spell checkers
    - Predictive text; auto-complete
    - Generative LLMs

### Group discussion

- *(Put into the shared deck)* Come up with at least one example
- How often do you use tools to check the spelling errors?

- English spelling remains complex
- Thankfully, we have **writers' aids**:
    - Spell checkers
    - Predictive text; auto-complete
    - Generative LLMs

### Group discussion

- *(Put into the shared deck)* Come up with at least one example
- How often do you use tools to check the spelling errors?
- Which one do you rely on the most?

# Spelling in the digital age

- English spelling remains complex
- Thankfully, we have **writers' aids**:
  - Spell checkers
  - Predictive text; auto-complete
  - Generative LLMs

**Group discussion**

- *(Put into the shared deck)* Come up with at least one example
- How often do you use tools to check the spelling errors?
- Which one do you rely on the most?
- Do they ever create problems (instead of solving them)?

Not all spelling errors are the same.

Not all spelling errors are the same.
To solve them, we need to consider **error types**.

# Different types of spelling errors

- 1. Non-word errors
- 2. Real-word errors
- *Notes.* How common?

· True confusion:

- True confusion:



- *"sissors"* (not knowing the correct form)

- True confusion:



- "*sissors*" (not knowing the correct form)
- **Typos**: "*hte*" (keyboard slip)

- True confusion:



- "*sissors*" (not knowing the correct form)
- **Typos**: "*hte*" (keyboard slip)
- Automatically detected when:

- True confusion:



- "*sissors*" (not knowing the correct form)
- **Typos**: "*hte*" (keyboard slip)
- Automatically detected when:
  - Word *not found* in **dictionary** of correct spellings

- Measures how "far apart" two strings are

- Measures how "far apart" two strings are
- Known as **Levenshtein distance**

- Measures how "far apart" two strings are
- Known as **Levenshtein distance**
- Minimum number of operations to transform one word into another

# Edit distance: Basic operations

Each operation = 1 unit of cost

- **Insertion**: *aquire → ac_c_quire*
- **Deletion**: *argu_e_ment → argument*
- **Substitution**: *cal_e_nder → cal_a_ndar*
- **Transposition**: *con_cs_ious → con_sc_ious*
    - Sometimes counted as two substitutions

- Helps suggest the **closest correct word** when a typo is found

- Helps suggest the **closest correct word** when a typo is found
- In other words, can be used to suggest candidate corrections

- Helps suggest the **closest correct word** when a typo is found
- In other words, can be used to suggest candidate corrections
  1. Input: *recieve*

- Helps suggest the **closest correct word** when a typo is found
- In other words, can be used to suggest candidate corrections
  1. Input: *recieve*
     - Candidates: *receive, recipe*

- Helps suggest the **closest correct word** when a typo is found
- In other words, can be used to suggest candidate corrections
  1. Input: *recieve*
     - Candidates: *receive, recipe*
  2. Input: *acommodation*

- Helps suggest the **closest correct word** when a typo is found
- In other words, can be used to suggest candidate corrections
  1. Input: *recieve*
     - Candidates: *receive, recipe*
  2. Input: *acommodation*
     - Candidates: *accommodation, commendation*

- Not all errors are equally likely

- Not all errors are equally likely
- Edit distance can be **weighted** for more realistic corrections

- Not all errors are equally likely
- Edit distance can be **weighted** for more realistic corrections
- Substituting a nearby key on the keyboard may cost less than a distant one

- Not all errors are equally likely
- Edit distance can be **weighted** for more realistic corrections
- Substituting a nearby key on the keyboard may cost less than a distant one
  - e.g., *friemd → friend* (substitution: *m→n*, keys are adjacent → low cost)

- Not all errors are equally likely
- Edit distance can be **weighted** for more realistic corrections
- Substituting a nearby key on the keyboard may cost less than a distant one
  - e.g., *friemd → friend* (substitution: *m→n*, keys are adjacent → low cost)
  - vs. *friemd → fried* (deletion of *m*, more disruptive → higher cost)

**Traditional method**: Dictionary + Edit Distance: How it works

- Relies on a dictionary of correct words

Limitations

**Traditional method**: Dictionary + Edit Distance: How it works

- Relies on a dictionary of correct words
- Calculates distance between misspelling and candidates

Limitations

**Traditional method**: Dictionary + Edit Distance: How it works

- Relies on a dictionary of correct words
- Calculates distance between misspelling and candidates
- Suggests the closest candidate as the correction

Limitations

**Traditional method**: Dictionary + Edit Distance: How it works

- Relies on a dictionary of correct words
- Calculates distance between misspelling and candidates
- Suggests the closest candidate as the correction
- Adds some weights for more realistic correction

**Limitations**

# Interim summary

**Traditional method**: Dictionary + Edit Distance: How it works

- Relies on a dictionary of correct words
- Calculates distance between misspelling and candidates
- Suggests the closest candidate as the correction
- Adds some weights for more realistic correction

## Limitations

- Fails with new words or domain-specific terms (e.g., *rizz*, *COVID-19*)

**Traditional method**: Dictionary + Edit Distance: How it works

- Relies on a dictionary of correct words
- Calculates distance between misspelling and candidates
- Suggests the closest candidate as the correction
- Adds some weights for more realistic correction

Limitations

- Fails with new words or domain-specific terms (e.g., *rizz*, *COVID-19*)
- Ignores context (e.g., *I want to by a book* → intended: *buy*)

**Traditional method**: Dictionary + Edit Distance: How it works

- Relies on a dictionary of correct words
- Calculates distance between misspelling and candidates
- Suggests the closest candidate as the correction
- Adds some weights for more realistic correction

### Limitations

- Fails with new words or domain-specific terms (e.g., *rizz*, *COVID-19*)
- Ignores context (e.g., *I want to by a book* → intended: *buy*)

**Traditional method**: Dictionary + Edit Distance: How it works

- Relies on a dictionary of correct words
- Calculates distance between misspelling and candidates
- Suggests the closest candidate as the correction
- Adds some weights for more realistic correction

Limitations

- Fails with new words or domain-specific terms (e.g., *rizz*, *COVID-19*)
- Ignores context (e.g., *I want to by a book* → intended: *buy*)

Q. What happens if the misspelled word is still a real word?

1. Local syntactic errors: *<u>Their</u> was a problem*

1. Local syntactic errors: *_Their_ was a problem*
2. Long-distance syntactic errors: *The key to the cabinets _are_ on the table*

1. Local syntactic errors: *_Their_ was a problem*
2. Long-distance syntactic errors: *The key to the cabinets _are_ on the table*
3. Semantic errors: *I read the _brook_*

1. Local syntactic errors: *<u>Their</u> was a problem*
2. Long-distance syntactic errors: *The key to the cabinets <u>are</u> on the table*
3. Semantic errors: *I read the <u>brook</u>*

- Solving this problem is more difficult:

1. Local syntactic errors: *<u>Their</u> was a problem*
2. Long-distance syntactic errors: *The key to the cabinets <u>are</u> on the table*
3. Semantic errors: *I read the <u>brook</u>*

- Solving this problem is more difficult:
  - The result is still a valid word → not flagged by a dictionary

## 2. Real-word Errors

1. Local syntactic errors: *Their was a problem*
2. Long-distance syntactic errors: *The key to the cabinets are on the table*
3. Semantic errors: *I read the brook*

- Solving this problem is more difficult:
    - The result is still a valid word → not flagged by a dictionary
    - Surrounding **context** must be considered

# How common are spelling errors?

- About 2–3% of all typed words on a full-size keyboard are misspelled by proficient adults (Flor et al., 2015)

Table 2. Summary statistics for the ETS Spelling Corpus

| | GRE Argument | GRE Issue | TOEFL Independent | TOEFL Integrated | TOTAL |
|---|---|---|---|---|---|
| Total essays | 750 | 750 | 750 | 750 | 3,000 |
| Essays without misspellings | 60 | 21 | 18 | 21 | 120 |
| Total Word Count | 263,578 | 336,301 | 212,930 | 151,031 | 963,840 |
| Average Word Count | 351 | 448 | 284 | 201 | 321 |
| Total count of Misspellings | 5,935 | 7,962 | 7,285 | 5,230 | 26,412 |
| Misspellings as % of all words | 2.25% | 2.37% | 3.42% | 3.46% | 2.74% |

**Figure 1:** Flor et al. (2015), p. 112

- About 2–3% of all typed words on a full-size keyboard are misspelled by proficient adults (Flor et al., 2015)

*Table 2. Summary statistics for the ETS Spelling Corpus*

|  | GRE Argument | GRE Issue | TOEFL Independent | TOEFL Integrated | TOTAL |
|---|---|---|---|---|---|
| Total essays | 750 | 750 | 750 | 750 | 3,000 |
| Essays without misspellings | 60 | 21 | 18 | 21 | 120 |
| Total Word Count | 263,578 | 336,301 | 212,930 | 151,031 | 963,840 |
| Average Word Count | 351 | 448 | 284 | 201 | 321 |
| Total count of Misspellings | 5,935 | 7,962 | 7,285 | 5,230 | 26,412 |
| Misspellings as % of all words | 2.25% | 2.37% | 3.42% | 3.46% | 2.74% |

**Figure 1:** Flor et al. (2015), p. 112

- Most errors are single-character misspellings (edit distance = 1)

- About 2–3% of all typed words on a full-size keyboard are misspelled by proficient adults (Flor et al., 2015)

*Table 2. Summary statistics for the ETS Spelling Corpus*

| | GRE Argument | GRE Issue | TOEFL Independent | TOEFL Integrated | TOTAL |
|---|---|---|---|---|---|
| Total essays | 750 | 750 | 750 | 750 | 3,000 |
| Essays without misspellings | 60 | 21 | 18 | 21 | 120 |
| Total Word Count | 263,578 | 336,301 | 212,930 | 151,031 | 963,840 |
| Average Word Count | 351 | 448 | 284 | 201 | 321 |
| Total count of Misspellings | 5,935 | 7,962 | 7,285 | 5,230 | 26,412 |
| Misspellings as % of all words | 2.25% | 2.37% | 3.42% | 3.46% | 2.74% |

**Figure 1:** Flor et al. (2015), p. 112

- Most errors are single-character misspellings (edit distance = 1)
- On a mobile phone, however, about 40% of words are misspelled (Grammarly, 2019)

- About 2–3% of all typed words on a full-size keyboard are misspelled by proficient adults (Flor et al., 2015)

Table 2. Summary statistics for the ETS Spelling Corpus

| | GRE Argument | GRE Issue | TOEFL Independent | TOEFL Integrated | TOTAL |
|---|---|---|---|---|---|
| Total essays | 750 | 750 | 750 | 750 | 3,000 |
| Essays without misspellings | 60 | 21 | 18 | 21 | 120 |
| Total Word Count | 263,578 | 336,301 | 212,930 | 151,031 | 963,840 |
| Average Word Count | 351 | 448 | 284 | 201 | 321 |
| Total count of Misspellings | 5,935 | 7,962 | 7,285 | 5,230 | 26,412 |
| Misspellings as % of all words | 2.25% | 2.37% | 3.42% | 3.46% | 2.74% |

**Figure 1:** Flor et al. (2015), p. 112

- Most errors are single-character misspellings (edit distance = 1)
- On a mobile phone, however, about 40% of words are misspelled (Grammarly, 2019)
- More multi-error misspellings and real-word errors due to auto-complete (e.g., *restaurant* → typed as *restuarnt* → auto-corrected to *restart*)

21

# Building a simple spell-checker

# Baseline spell checker (Peter Norvig)

- Generate all candidate words within 1–2 edits

# Baseline spell checker (Peter Norvig)

- Generate all candidate words within 1–2 edits
- Keep only words in the *dictionary* from a *corpus*

## Baseline spell checker (Peter Norvig)

- Generate all candidate words within 1–2 edits
- Keep only words in the *dictionary* from a *corpus*
- Pick the most frequent candidate

## Baseline spell checker (Peter Norvig)

- Generate all candidate words within 1–2 edits
- Keep only words in the *dictionary* from a *corpus*
- Pick the most frequent candidate
- Example:

- Generate all candidate words within 1–2 edits
- Keep only words in the *dictionary* from a *corpus*
- Pick the most frequent candidate

- Example:
  - Input: *langage*

## Baseline spell checker (Peter Norvig)

- Generate all candidate words within 1–2 edits
- Keep only words in the *dictionary* from a *corpus*
- Pick the most frequent candidate

- Example:
    - Input: *langage*
    - Candidates: *language*, *lineage*

## Baseline spell checker (Peter Norvig)

- Generate all candidate words within 1–2 edits
- Keep only words in the *dictionary* from a *corpus*
- Pick the most frequent candidate

- Example:
  - Input: *langage*
  - Candidates: *language*, *lineage*
  - Output: language

- Some typos are more likely than others

- Some typos are more likely than others
  - Adjacent key slips (e.g., *friemd* → *friend*)

# But there's a problem...

- Some typos are more likely than others
    - Adjacent key slips (e.g., *friemd* → *friend*)
    - Transpositions (e.g., *teh* → *the*)

## But there's a problem…

- Some typos are more likely than others
  - Adjacent key slips (e.g., *friemd* → *friend*)
  - Transpositions (e.g., *teh* → *the*)
- Baseline only looks at **frequency**, not how errors happen

# But there's a problem...

- Some typos are more likely than others
    - Adjacent key slips (e.g., *friemd* → *friend*)
    - Transpositions (e.g., *teh* → *the*)
- Baseline only looks at **frequency**, not how errors happen
- We need a better model: **noisy channel**

Formula
$\arg \max_w P(\text{observed} \mid w) \cdot P(w)$

Formula
$\arg\max_w P(\text{observed} \mid w) \cdot P(w)$

- P(w) = prior probability (word frequency)

# Noisy channel spell checker

Formula
$\arg \max_w P(\text{observed} \mid w) \cdot P(w)$

- **P(w)** = prior probability (word frequency)
- **P(observed|w)** = likelihood of making that typo

#### Formula
arg max$_w$ $P(\text{observed} \mid w) \cdot P(w)$

- **P(w)** = prior probability (word frequency)
- **P(observed|w)** = likelihood of making that typo
- Example:

# Noisy channel spell checker

**Formula**
$\arg\max_w P(\text{observed} \mid w) \cdot P(w)$

- **P(w)** = prior probability (word frequency)
- **P(observed|w)** = likelihood of making that typo
- Example:
    - Input: *recieve*

# Noisy channel spell checker

### Formula
$\arg\max_w P(\text{observed} \mid w) \cdot P(w)$

- **P(w)** = prior probability (word frequency)
- **P(observed|w)** = likelihood of making that typo
- Example:
    - Input: *recieve*
    - Candidates: *recipe*, receive

# Noisy channel spell checker

### Formula
$\arg \max_w P(\text{observed} \mid w) \cdot P(w)$

- **P(w)** = prior probability (word frequency)
- **P(observed|w)** = likelihood of making that typo
- Example:
    - Input: *recieve*
    - Candidates: *recipe*, **receive**
    - Baseline (frequency only) → *recipe*

# Noisy channel spell checker

### Formula
arg max$_w$ $P(\text{observed} \mid w) \cdot P(w)$

- **P(w)** = prior probability (word frequency)
- **P(observed|w)** = likelihood of making that typo
- Example:
    - Input: *recieve*
    - Candidates: *recipe*, **receive**
    - Baseline (frequency only) → *recipe*
    - Noisy channel (frequency + typo likelihood) → **receive**

# Thinking about a more complex spell-checker

Example: Someone types:

```
You put the catt before the horse.
```

· The word `catt` is not found in a dictionary or corpus ⇒ likely a misspelling.

## Why context matters in spell-checking

Example: Someone types:

    You put the `catt` before the horse.

- The word `catt` is not found in a dictionary or corpus ⇒ likely a misspelling.
- A simple spell-checker (like Norvig's) would:

Example: Someone types:

You put the catt before the horse.

- The word catt is not found in a dictionary or corpus $\Rightarrow$ likely a misspelling.
- A simple spell-checker (like Norvig's) would:
  - Consider all known words with edit distance 1 from catt

## Why context matters in spell-checking

Example: Someone types:

> You put the `catt` before the horse.

- The word `catt` is not found in a dictionary or corpus $\Rightarrow$ likely a misspelling.
- A simple spell-checker (like Norvig's) would:
  - Consider all known words with edit distance 1 from `catt`
  - Choose the most frequent candidate: `cat`

**Example:** Someone types:

    You put the catt before the horse.

- The word `catt` is not found in a dictionary or corpus ⇒ likely a misspelling.
- A simple spell-checker (like Norvig's) would:
    - Consider all known words with edit distance 1 from `catt`
    - Choose the most frequent candidate: `cat`
- But the better correction is actually `cart`, because:

**Example:** Someone types:

`You put the ` `catt` ` before the horse.`

- The word `catt` is not found in a dictionary or corpus $\Rightarrow$ likely a misspelling.
- A simple spell-checker (like Norvig's) would:
  - Consider all known words with edit distance 1 from `catt`
  - Choose the most frequent candidate: `cat`
- But the better correction is actually `cart`, because:
  - `put the ` `cart` ` before the horse` is a common English expression

**Example:** Someone types:

```
You put the catt before the horse.
```

- The word `catt` is not found in a dictionary or corpus ⇒ likely a misspelling.
- A simple spell-checker (like Norvig's) would:
  - Consider all known words with edit distance 1 from `catt`
  - Choose the most frequent candidate: `cat`
- But the better correction is actually `cart`, because:
  - `put the cart before the horse` is a common English expression
  - `put the cat before the horse` is not

N-grams are sequences of *n* elements (e.g., words or characters):

· Unigram = one word: `the`

## Using N-grams to model context

N-grams are sequences of *n* elements (e.g., words or characters):

- **Unigram** = one word: `the`
- **Bigram** = two-word sequence: `the cat`

## Using N-grams to model context

N-grams are sequences of *n* elements (e.g., words or characters):

- **Unigram** = one word: `the`
- **Bigram** = two-word sequence: `the cat`
- **Trigram** = three-word sequence: `put the cat`

## Using N-grams to model context

N-grams are sequences of *n* elements (e.g., words or characters):

- **Unigram** = one word: `the`
- **Bigram** = two-word sequence: `the cat`
- **Trigram** = three-word sequence: `put the cat`

N-grams are sequences of *n* elements (e.g., words or characters):

- **Unigram** = one word: `the`
- **Bigram** = two-word sequence: `the cat`
- **Trigram** = three-word sequence: `put the cat`

#### How do we use this?

- Count all *n*-grams (e.g., bigrams) in a large corpus

Demo: `https://huggingface.co/spaces/liujch1998/infini-gram`

## Using N-grams to model context

N-grams are sequences of *n* elements (e.g., words or characters):

- **Unigram** = one word: `the`
- **Bigram** = two-word sequence: `the cat`
- **Trigram** = three-word sequence: `put the cat`

#### How do we use this?

- Count all *n*-grams (e.g., bigrams) in a large corpus
- Use frequency of phrases to estimate how likely a candidate is *in context*

Demo: `https: //huggingface.co/spaces/liujch1998/infini-gram`

## Using N-grams to model context

N-grams are sequences of *n* elements (e.g., words or characters):

- **Unigram** = one word: `the`
- **Bigram** = two-word sequence: `the cat`
- **Trigram** = three-word sequence: `put the cat`

### How do we use this?

- Count all *n*-grams (e.g., bigrams) in a large corpus
- Use frequency of phrases to estimate how likely a candidate is *in context*
- `put the cart` is more frequent than `put the cat`

Demo: `https: //huggingface.co/spaces/liujch1998/infini-gram`

- **Statistical Language Models (n-grams)** Use probability of surrounding context e.g., *I went to the shcool* → "school" is more probable

**Summary:** Traditional = simple but context-blind
Modern = complex but context-aware

## Other approaches

- **Statistical Language Models (n-grams)** Use probability of surrounding context e.g., *I went to the shcool* → "school" is more probable
- **Neural Spell Checkers (Deep Learning)** Seq2Seq / Transformer-based models generate corrected text Examples: ChatGPT, Grammarly, Google Docs

**Summary:** Traditional = simple but context-blind
Modern = complex but context-aware

- **Statistical Language Models (n-grams)** Use probability of surrounding context e.g., *I went to the shcool* → "school" is more probable
- **Neural Spell Checkers (Deep Learning)** Seq2Seq / Transformer-based models generate corrected text Examples: ChatGPT, Grammarly, Google Docs
- **Hybrid Approaches** Combine edit distance with language models; pick the highest probability candidate

**Summary:** Traditional = simple but context-blind
Modern = complex but context-aware

## Real-Word Errors and Grammar

- Not all mistakes are spelling errors → some are **real-word errors**.

# Real-Word Errors and Grammar

- Not all mistakes are spelling errors → some are **real-word errors**.
  - Example: *I want to by a book* → "by" is valid, intended: *buy*

## Real-Word Errors and Grammar

- Not all mistakes are spelling errors → some are **real-word errors**.
  - Example: *I want to by a book* → "by" is valid, intended: *buy*
- Real-word errors often overlap with **grammar errors**.

# Real-Word Errors and Grammar

- Not all mistakes are spelling errors → some are **real-word errors**.
    - Example: *I want to by a book* → "by" is valid, intended: *buy*
- Real-word errors often overlap with **grammar errors**.
    - Example: *Their going to school* → all words exist, but grammar is wrong (*They're*)

## Real-Word Errors and Grammar

- Not all mistakes are spelling errors → some are **real-word errors**.
    - Example: *I want to by a book* → "by" is valid, intended: *buy*
- Real-word errors often overlap with **grammar errors**.
    - Example: *Their going to school* → all words exist, but grammar is wrong (*They're*)
- Modern systems therefore blur the line between spell checking and grammar checking, using **context-aware models** to handle both (which we'll talk about in the next class).

# Wrap-up

## Wrap-up

Key idea: ~~Spelling errors are annoying~~

# Wrap-up

Key idea: ~~Spelling errors are annoying~~
Spelling errors vary by types

Key idea: ~~Spelling errors are annoying~~
Spelling errors vary by types More questions to think about:

- What about the spacing errors?
- What about in other languages that have different encoding systems?