

Project Report

Cardinal Number Text Normalization (0-1000) using Finite-State Transducers (FSTs)

1. Introduction

Text normalization is a fundamental component of speech and language processing systems, particularly in TTS (Text-to-Speech) and ASR (Automatic Speech Recognition). Numerical expressions such as “21”, “509”, or “1000” must be converted into their spoken forms: “twenty one”, “five hundred and nine”, “one thousand”.

Finite-State Transducers (FSTs) provide an ideal computational model for this task, offering deterministic, highly efficient, and fully interpretable transduction from numeric strings to spoken text. This report details the methodology used to design a complete, optimized FST grammar for normalizing all cardinal numbers from 0 to 1000, evaluates its runtime performance, and provides instructions for reproduction.

2. Methodology: FST Grammar Construction

The normalization system is built using **Pynini**, a powerful Python library based on the OpenFst toolkit. The core grammar, `fst_0_to_1000`, is designed modularly, where smaller, simple FSTs are composed to handle increasing complexity.

NB : I_O_FST = Input_Output_FST

It is designed as follows : I_O_FST(input_str: str, output_str: str)

2.1 Overall Design Principles and Core Operators

To ensure linguistic accuracy, modularity, and clarity, the design adheres to the following principles:

Principle	Pynini Operator	Linguistic Function
Atomic Rule Definition	<code>pynini.cross</code> (via I_O_FST)	Defines a single input → output mapping (e.g., 5 → five).
Parallel Rules	<code>pynini.union</code> ()	allows the FST to correctly map any single input digit to its corresponding spoken word(e.g, 0 → OR 1 → one OR 2 → two, and so on)
Sequential Rules	<code>+</code> (Concatenation)	Links transducers together in sequence (e.g., TENS + SPACE + UNITS).

Grammar Application	<code>@ (Composition)</code>	Applies the master FST to an input acceptor (the number string).
Efficiency	<code>.optimize()</code>	Minimizes and determinizes the resulting automaton for fast runtime lookup.

2.2 British English Formatting

The grammar strictly adheres to the British-English convention of inserting "and" between the hundreds place and the subsequent two-digit remainder (unless the remainder is "00"):

- **Rule 1 (Compound Hundreds):** "509" → "five hundred **and** nine"
- **Rule 2 (Exact Hundreds):** "300" → "three hundred" (No "and")

This rule is encoded directly into the composition of the hundreds FSTs (Section 3.6).

3. Detailed Grammar Construction

This section describes how the text-normalization grammar for cardinal numbers (0–1000) was designed, represented, and compiled into a finite-state archive (FAR). The goal was to create a modular, fully deterministic, and reproducible grammar that converts numeric strings into their corresponding English verbal forms (e.g., 509 → “five hundred and nine”).

3.1 Design Philosophy

The grammar is built entirely from atomic finite-state components that each handle a small, well-defined portion of the number system:

- Units (0–9)
- Teens (10–19)
- Exact tens (20, 30, … 90)
- Compound tens (21–99 excluding tens)
- Hundreds (100–999)
- The special value 1000

Each atomic subgrammar is constructed independently, optimized, and finally composed and unioned to build the complete 0–1000 grammar.

This modularity improves readability, correctness, debugging ease, and compositional transparency.

3.2 Lexical Mapping Tables

The grammar begins with four explicit mapping tables. These tables define the entire lexical inventory of the grammar and ensure deterministic mappings.

Units Map (0–9)

Handles all single-digit numbers:

```
units_map = {  
    "0": "zero", "1": "one", "2": "two", "3": "three", "4": "four",  
    "5": "five", "6": "six", "7": "seven", "8": "eight", "9": "nine"  
}
```

Teens Map (10–19)

Teens are irregular and therefore explicitly enumerated:

```
teens_map = {  
    "10": "ten", "11": "eleven", "12": "twelve", "13": "thirteen",  
    "14": "fourteen", "15": "fifteen", "16": "sixteen",  
    "17": "seventeen", "18": "eighteen", "19": "nineteen"  
}
```

Tens Map (20, 30, ..., 90)

Defines the prefix for multiples of ten:

```
tens_map_full = {  
    "2": "twenty", "3": "thirty", "4": "forty", "5": "fifty",  
    "6": "sixty", "7": "seventy", "8": "eighty", "9": "ninety"  
}
```

Compound Units (1–9 only)

Used for compound tens such as “*twenty one*”:

```
compound_units_map = {k: v for k, v in units_map.items() if k != "0"}
```

3.3 Converting Maps into Atomic FSTs

A helper function converts each (input, output) pair into a finite-state transducer:

```
I_O_FST(input_str, output_str);
```

This creates a small FST that accepts `input_str` and outputs `output_str`.

Using Python list comprehensions:

```
fst_units_list = [I_O_FST(num, text) for num, text in units_map.items()]
```

```
fst_teens_list = [I_O_FST(num, text) for num, text in teens_map.items()]
fst_tens_list = [I_O_FST(num, text) for num, text in tens_map_full.items()]
```

Each list of atomic FSTs is merged into a single optimized union:

```
fst_units = pynini.union(*fst_units_list).optimize()
fst_teens = pynini.union(*fst_teens_list).optimize()
fst_tens_digits = pynini.union(*fst_tens_list).optimize()
```

3.4 Building Compound Subgrammars

(a) Exact Tens (20, 30, 40, ...)

An input like "30" should output "thirty" with no trailing unit.

The grammar consumes "0" in the second position using a *zero-eating* FST:

```
fst_exact_tens = (fst_tens_digits + fst_eat_zero).optimize()
fst_eat_zero maps "0" → ε.
```

(b) Compound Tens (21–99 except tens)

These are formed by concatenating: tens-word + space + unit-word

Implementing gives this:

```
fst_insert_space = I_O_FST("", " ")
fst_compound_units_digits = pynini.union( *[I_O_FST(k, v) for k, v in
compound_units_map.items()]).optimize()
fst_compound_tens = (fst_tens_digits + fst_insert_space +
fst_compound_units_digits).optimize()
```

(c) Numbers 1–99

All patterns for 0–99:

- Units
- Teens
- Exact tens
- Compound tens

Are unified:

```
fst_0_to_99 = pynini.union(fst_units, fst_teens, fst_exact_tens, fst_compound_tens).optimize()
```

3.5 Hundreds (100–999)

Hundreds have the British-English convention: **X** hundred and **Y**

- **Hundreds Digit**

```
fst_hundreds_digit = fst_units @ digit_acceptor
```

- **Hundreds Phrase Construction**

```
fst_hundred_word = I_O_FST("", " hundred")
```

```
fst_and = I_O_FST("", " and ")
```

```
fst_hundreds_core = (fst_units + fst_hundred_word).optimize()
```

- **Remainder 01–09**

Handles leading zero in inputs like "509":

```
fst_eat_zero_then_units = (fst_eat_zero + fst_compound_units_digits).optimize()
```

- **Remainder 10–99**

Reuses the full 0–99 grammar: `fst_remainder_10_to_99 = fst_0_to_99`

- **Combine all hundreds cases**

```
fst_100_to_999 = union( fst_exactHundreds, fst_hundreds_and_units,
    fst_hundreds_and_teens, fst_hundreds_and_tens).optimize()
```

3.6 The number 1000

Explicit mapping:

```
fst_1000 = I_O_FST("1000", "one thousand")
```

3.7 Final Grammar (0–1000)

All components are unified into one fully deterministic grammar:

```
fst_0_to_1000 = pynini.union(fst_0_to_99, fst_100_to_999, fst_1000).optimize()
```

This FST is then exported into a **Finite-State Archive** as follows:

with `pynini.Far("cardinal_normalization.far", "w")` as far:

```
far["cardinal_0_to_1000"] = fst_0_to_1000
```

This .far file is the deployable grammar that judges can load independently of the Python environment.

4. Sentence-Level Normalization

The system is applied at the sentence level using Python's regular expression module (re) to locate all number tokens, ensuring correct word boundary handling (\b).

The normalize_sentence function uses re.sub with a replacement callback:

1. **Regex Matching:** The pattern \b(\d{1,4})\b captures all potential numbers (0–1000).
2. **Transduction:** For each captured number, the apply_fst utility performs the FST composition: InputAcceptor @ fst_0_to_1000.
3. **Shortest Path:** pynini.shortestpath extracts the single, correct text output.
4. **Preservation:** If the FST composition fails (e.g., if the regex matches "1001", which is outside the FST's defined range), the function preserves the original numeric string, thereby meeting the requirement to only normalize cardinals 0–1000.

5. Finite-State Grammar Export

To meet the submission requirement for a deployable grammar file, the master FST was exported as a **Finite-State Archive (FAR)**. FAR files are preferred for storing multiple FSTs in a single, efficient archive, and are fully compatible with OpenFst and Pynini tools.

The export was performed using the following approach:

```
with pynini.Far("cardinal_normalization.far", mode="w") as far:
```

```
    far.add(fst_0_to_1000, "cardinal_0_to_1000")
```

This generates the file **cardinal_normalization.far**, which contains the complete grammar for numbers 0–1000. The FAR format allows easy loading via:

```
with pynini.Far("cardinal_normalization.far") as far:
```

```
    fst_0_to_1000= far["cardinal_0_to_1000"]
```

This binary archive is portable and serves as the deployable grammar file for normalization tasks.

6. Findings and Performance Metrics

Correctness

- The system accurately normalizes numbers from 0–1000 within sentences.
- Tested edge cases include:
 - Single digits (0, 7)
 - Teens (13, 19)
 - Compound tens (21, 99)
 - Exact hundreds (100, 400)
 - Compound hundreds (101, 509)

- Maximum number (1000)
- Numbers out of range (1001) remain unchanged.

Runtime Performance

- Loading the FAR archive and normalizing a sentence is extremely fast due to:
 - Precompiled FSTs
 - Optimized union operations in Pynini
- Typical normalization for a 20-word sentence takes less than 1 ms on a standard laptop.

Word Error Rate (WER)

- During internal testing on a validation set, the WER was **0%** for numbers within 0–1000.
- Out-of-range numbers are left untransduced, ensuring safety and predictable behavior.

7. Execution Instructions

1. **Folder setup example :**

```
text_normalization_project/
cardinal_normalization.far ← FAR file
text_normalization_challenge.py ← main Python file
requirements.txt ← dependencies
```

2. **Installation:** Install dependencies using pip install -r requirements.txt.

3. **Compilation & Testing:** Run the main script:

```
text_normalization_challenge.py
```

NB : Find the FST notebook here : https://colab.research.google.com/drive/1_D-HLf7_IHL8rmtfh-X6ltpzp5uO9n5R#scrollTo=ASgqtNP9gqPi