# Non-Regular Languages

Monday September 27, 2021

Turing-recognizable

decidable

context-free

regular

???

# Announcements

- HW2 due yesterday

- HW3 released, due Sun 10/3 11:59pm EST

- First in-person class: next Monday 10/4
  - McCormack M01-0209

# *So Far:* Regular or Not?

- Many ways to prove that a language is regular:
    - Construct a <u>DFA</u> or <u>NFA</u> (or GNFA)
    - Come up with a <u>regular expression</u> describing the language

- But how to show that a language is **not regular**?
    - E.g., HTML / XML is <u>not</u> a regular language
    - Can't be represented with a regular expression (common mistake)!

# *Flashback:* Designing DFAs or NFAs

- Each state "stores" some information
  - E.g., $q_0$ = "seen zero **1s**", $q_1$ = "seen one **1**", $q_2$ = "seen two **1s**" etc.
  - <u>Finite</u> states = <u>finite</u> amount of info (decided in advance)

- This means <u>DFAs can't keep track of an arbitrary count</u>!
  - would require infinite states

169

# A Non-Regular Language

$L = \{\ \mathbf{0}^n\mathbf{1}^n\ |\ n >= \mathbf{0}\ \}$

- A DFA recognizing $L$ would require infinite states! (impossible)
  - States representing zero $\mathbf{0}$s, one $\mathbf{0}$, two $\mathbf{0}$s, …

- This language represents the essence of many PLs, e.g., HTML!
  - To better see this replace:
    - "`0`" -> "`<tag>`" or "`(`"
    - "`1`" -> "`</tag>`" or "`)`"

- The problem is tracking the **<u>nestedness</u>**
  - Regular languages cannot count arbitrary nesting depths
  - So most programming language syntax is not regular!

Still, how do we
**prove non-regularness**?

# A Lemma About Regular Languages

**Pumping lemma**  If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0$, $xy^i z \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

All regular languages satisfy these three conditions!

Specifically, strings in the language longer than length $p$ satisfy the conditions

**Lemma doesn't tell you an exact $p$!**
(just that there exists "some" $p$)

# The Pumping Lemma: Finite Langu

**Pumping lemma** If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0$, $xy^i z \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

What's a possible $p$?
**Length of longest string + 1**

# strings in the language with at least length $p$? **None!**

Therefore, <u>all</u> strings with length at least $p$ satisfy the pumping lemma conditions! ☺

In finite langs, these are true for all strings "of length at least $p$" (for some $p$)

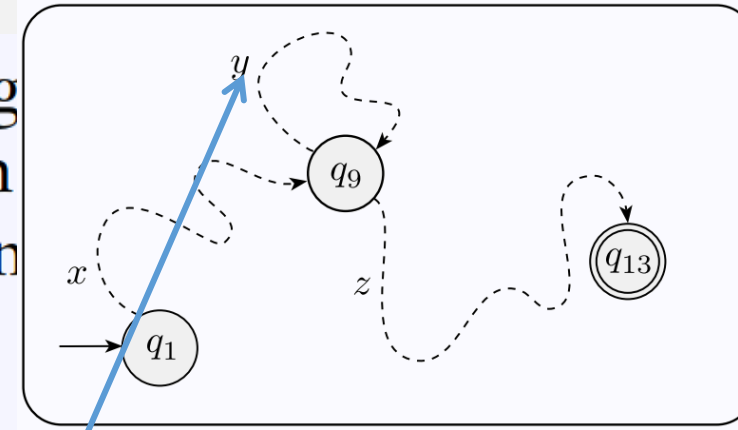Example: a finite language {"ab", "cd"}
- <u>All finite langs are regular</u> (can easily construct DFA/NFA recognizing them)

172

# The Pumping Lemma, a Closer Look

**Pumping lemma**     If $A$ is a regular lang[...]nber $p$ (the pumping length) where if $s$ is any string in [...] $s$ may be divided into three pieces, $s = xyz$, satisfyin[...]s:

1. for each $i \geq 0$, $xy^i z \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.



"long enough" strings, should have a __repeatable__ ("pumpable") part; "pumped" string is still in the language

## Strings that have a <u>repeatable</u> part can be split into:
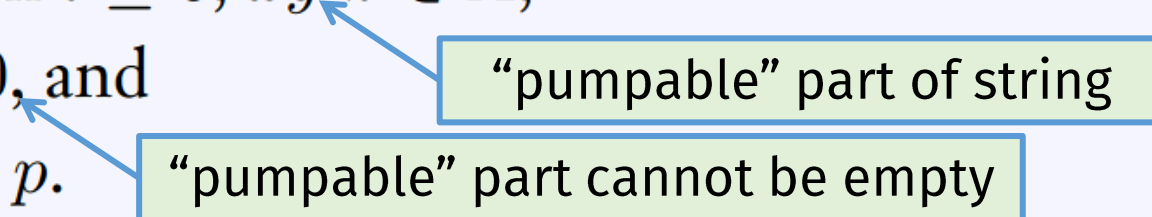
- $x$ = the part <u>before</u> any repeating
- $y$ = the repeated part
- $z$ = the part <u>after</u> any repeating

This makes sense because DFAs have a finite number of states, so for "long enough" (i.e., some length $p$) inputs, <u>some state must repeat</u>

e.g., "long enough length" = **# of states +1**
(The Pigeonhole Principle)

# The Pumping Lemma: Infinite Languages

**Pumping lemma**    If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

**1.** for each $i \geq 0$, $xy^i z \in A$,

**2.** $|y| > 0$, and

**3.** $|xy| \leq p$.

"pumpable" part of string

"pumpable" part cannot be empty

Example: *infinite* language {"00", "010" , "0110" , "01110", ...}
- Language is regular bc it's described by the regular expression 01*0
- Notice that the middle part is pumpable!
- E.g., "010" in the language can be split into three parts: x = 0, y = 1, z = 0
  - Any pumping (repeating) of the middle part creates a string that is still in the language
    - $i$ = 1 -> "0**1**0", $i$ = 2 -> "0**11**0", $i$ = 3 -> "0**111**0"

# Summary: The Pumping Lemma ...

- ... states properties that are <u>true for all regular languages</u>

**IMPORTANT:**
- The Pumping Lemma <u>cannot prove</u> that a language is regular!

- But ... we can use it to prove that a language is <u>not</u> regular

# Poll: Conditional Statements

# Equivalence of Conditional Statements

- Yes or No? "If X then Y" is equivalent to:

  - "If Y then X" (converse)
    - No!

  - "If not X then not Y" (inverse)
    - No!

  - "If not Y then not X" (contrapositive) ← Proof by contradiction
    - Yes!

184

# Pumping Lemma: Proving Non-Regularity

… then the language is **not** regular

**Pumping lemma**   If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0$, $xy^i z \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

If any of these are **not** true …

Contrapositive:
"If X then Y" is equivalent to "If **not** Y then **not** X"

# Pumping Lemma: Non-Regularity Example

Let $B$ be the language $\{0^n1^n \mid n \geq 0\}$. We use the pumping lemma to prove that $B$ is not regular. The proof is by contradiction.

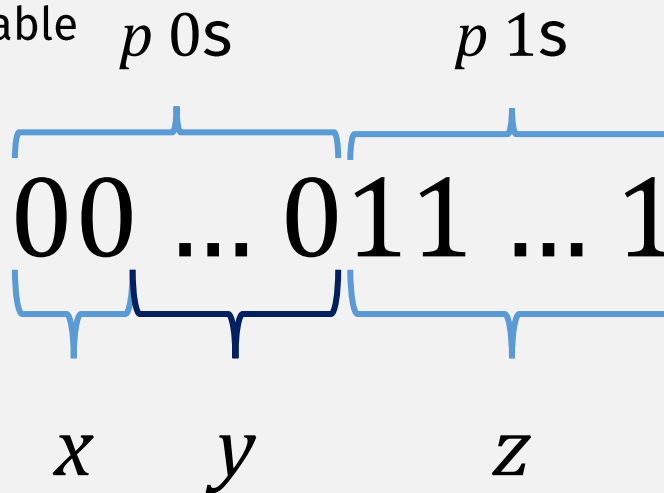# How To Do Proof By Contradiction

- Assume the opposite of the statement to prove

- Show that the assumption leads to a contradiction

- Conclude that the original statement must be true

# Possible Split: $y$ = all 0s

<u>Proof</u> (by contradiction):

- <u>Assume</u>: $0^n1^n$ **is** a regular language
  - So it must satisfy the pumping lemma
  - I.e., all strings length $p$ or longer are pumpable
- <u>Counterexample</u> = $0^p1^p$

- Choose $xyz$ split so $y$ contains:
  - all 0s

- Pumping $y$: produces a string with more 0s than 1s
  - Which is <u>not</u> in the language $0^n1^n$
  - This means that $0^p1^p$ does <u>not</u> satisfy the pumping lemma
  - Which means that that $0^n1^n$ is a <u>not</u> regular language
  - This is a **contradiction** of the assumption!

... then **not** true

**pumping lemma** → If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0$, $xy^iz \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

<u>Contrapositive</u>: If **not** true ...

<u>Reminder</u>: Pumping lemma says strings >= length $p$ splittable into $xyz$ where $y$ is pumpable

$p$ 0s          $p$ 1s

$$00 \ldots 011 \ldots 1$$

$x \quad y \quad\quad z$

<u>BUT</u> ... pumping lemma requires **only one** pumpable splitting

So the proof is not done!

Is there another way to split into $xyz$ ?

189

**Pumping lemma** If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0$, $xy^iz \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.
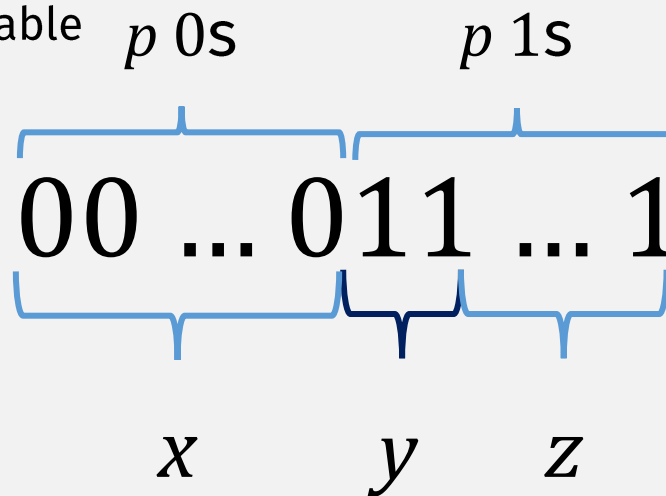
# Possible Split: $y$ = all 1s

Proof (by contradiction):

- <u>Assume</u>: $0^n1^n$ **is** a regular language
  - So it must satisfy the pumping lemma
  - I.e., all strings length $p$ or longer are pumpable
- <u>Counterexample</u> = $0^p1^p$

- Choose $xyz$ split so $y$ contains:
  - all 1s

$p$ 0s          $p$ 1s

$$00 \ldots 011 \ldots 1$$

$x \qquad y \qquad z$

- Is this string pumpable?
  - No!
  - By the same reasoning as in the previous slide

Is there another way to split into $xyz$ ?

190

# Possible Split: $y$ = 0s and 1s

**Pumping lemma** If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0$, $xy^i z \in A$,
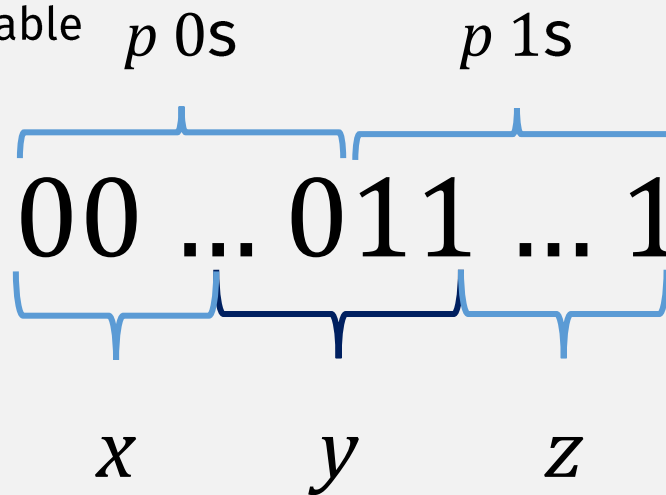2. $|y| > 0$, and
3. $|xy| \leq p$.

<u>Proof</u> (by contradiction):

- <u>Assume</u>: $0^n1^n$ **is** a regular language
  - So it must satisfy the pumping lemma
  - I.e., all strings length $p$ or longer are pumpable
- <u>Counterexample</u> = $0^p1^p$

$p$ 0s    $p$ 1s

- Choose $xyz$ split so $y$ contains:
  - both 0s and 1s

$$00 \ldots 011 \ldots 1$$

$x \qquad y \qquad z$

Did we examine every possible splitting?

**Yes! QED**

But maybe we did't have to ...

- Is this string pumpable?
  - No!
  - Pumped string will have equal 0s and 1s
  - But they will be in the wrong order: so there is still a **contradiction!**

191

# The Pumping Lemma: Condition 3

**Pumping lemma** If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

**1.** for each $i \geq 0$, $xy^iz \in A$,

**2.** $|y| > 0$, and

**3.** $|xy| \leq p$.

Repeating part $y$ ...
must be in the first $p$ characters!

$p$ 0s

$$00 \ldots 011 \ldots 1$$

$y$ must be in here!

# The Pumping Lemma: Pumping Down

**Pumping lemma**    If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0,$ $xy^i z \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

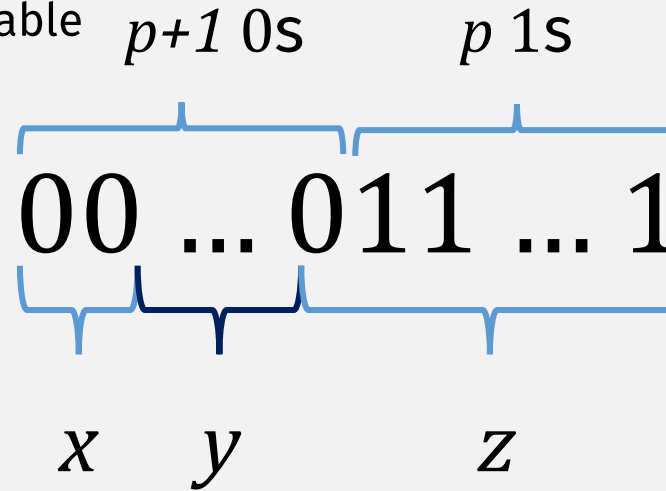Repeating part $y$ must be non-empty ...
but can be repeated zero times!

Example: $L = \{0^i 1^j \mid i > j\}$

# Pumping Down

Proof (by contradiction):

- Assume: $L$ **is** a regular language
  - So it must satisfy the pumping lemma
  - I.e., all strings length $p$ or longer are pumpable

- Counterexample = $0^{p+1}1^p$

$\overbrace{\qquad}^{\textit{p+1} \text{ 0s}} \overbrace{\qquad}^{\textit{p} \text{ 1s}}$

- Choose $xyz$ split so $y$ contains:
  - all 0s
  - (Only possibility, by condition 3)

$$00 \ldots 011 \ldots 1$$

$$x \quad y \quad\quad z$$

- Repeat $y$ zero times (**pump down**): produces string with 0s =< 1s
  - Which is <u>not</u> in the language $\{0^i 1^j \mid i > j\}$
  - This means that $\{0^i 1^j \mid i > j\}$ does <u>not</u> satisfy the pumping lemma
  - Which means that that it is a <u>not</u> regular language
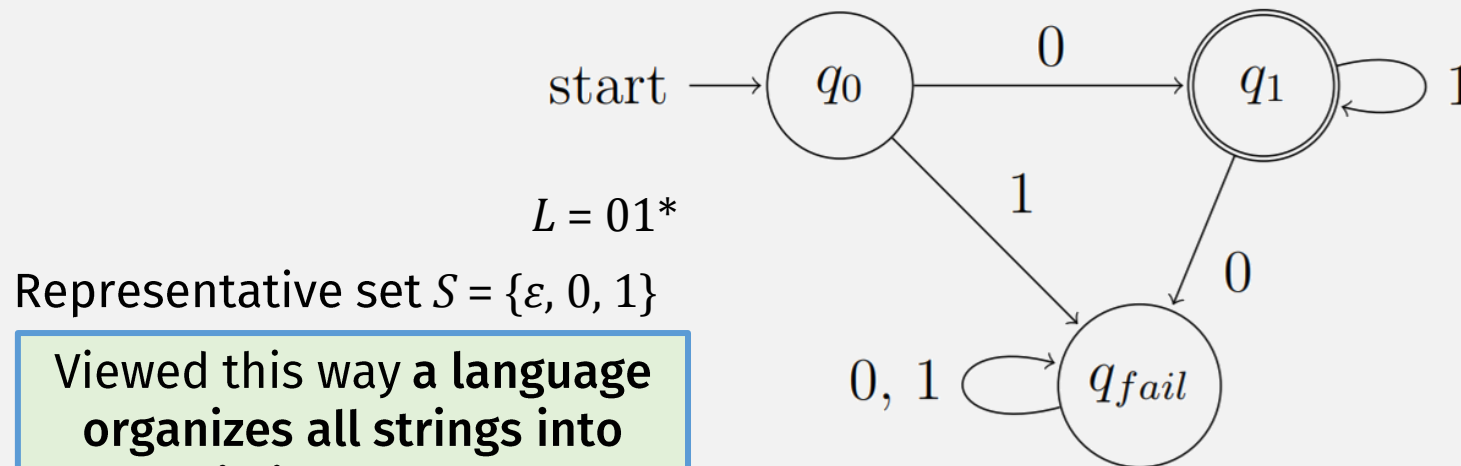  - This is a **contradiction** of the assumption!

# Pumping Lemma Doesn't Always Work!

- What if you can't figure out a counterexample?

# Another Way to Prove <u>Regularity</u>

- A set of strings $S$ is "representative" of a language $L$ if:
  - Every possible string $w \in \Sigma^*$ maps to a string $s$ in $S$ via REP where ...
  - $\text{REP}(w) = s$, if for every possible string $z$, $wz \in L$ iff $sz \in L$



$L = 01^*$

Representative set $S = \{\varepsilon, 0, 1\}$

For regular languages, strings in the "representative" set correspond to states in a DFA!

$S$ contains one string that reaches each state

Then $\text{REP}(w) = s$ if $w$ reaches the same state that $s$ represents

Viewed this way **a language organizes all strings into distinct groups**

Then for any string $z$, $wz \in L$ iff $sz \in L$ because they started in the same state!

**A language is regular if this number of groups is finite,** i.e. it has a finite representative set!

# Another Way to Prove <u>Non-Regularity</u>

- A set of strings $S$ is "representative" of a language $L$ if:
  - Every possible string $w \in \Sigma^*$ maps to a string $s$ in $S$ via REP where …
  - $\text{REP}(w) = s$, if for every possible string $z$, $wz \in L$ iff $sz \in L$

$L = \{ \mathbf{0}^n \mathbf{1}^n \mid n >= 0 \}$

- There must be a $\text{REP}(\mathbf{0}^k)$ every $k$ …
  - Because for every two strings $\mathbf{0}^k$ and $\mathbf{0}^m$ …
  - … there's some $z$ that completes it such that $\mathbf{0}^k z \in L$ but $\mathbf{0}^m z$ is not
  - E.g., let $z = \mathbf{1}^k$, then $\mathbf{0}^k \mathbf{1}^k \in L$ but $\mathbf{0}^m \mathbf{1}^k$ is not in $L$

The representative set is infinite!

So the language is not regular!

# Check-in Quiz 9/27

On gradescope