

Context-Free Grammars

Describing Languages

- We've seen two models for the regular languages:
 - ***Finite automata*** accept precisely the strings in the language.
 - ***Regular expressions*** describe precisely the strings in the language.
- Finite automata ***recognize*** strings in the language.
 - Perform a computation to determine whether a specific string is in the language.
- Regular expressions ***match*** strings in the language.
 - Describe the general shape of all strings in the language.

Context-Free Grammars

- A ***context-free grammar*** (or ***CFG***) is an entirely different formalism for defining a class of languages.
- ***Goal:*** Give a description of a language by recursively describing the structure of the strings in the language.
- CFGs are best explained by example...

Arithmetic Expressions

- Suppose we want to describe all legal arithmetic expressions using addition, subtraction, multiplication, and division.
- Here is one possible CFG:

$E \rightarrow \text{int}$
$E \rightarrow E \text{ Op } E$
$E \rightarrow (E)$
$\text{Op} \rightarrow +$
$\text{Op} \rightarrow -$
$\text{Op} \rightarrow *$
$\text{Op} \rightarrow /$

E
 $\Rightarrow E \text{ Op } E$
 $\Rightarrow E \text{ Op } (E)$
 $\Rightarrow E \text{ Op } (E \text{ Op } E)$
 $\Rightarrow E * (E \text{ Op } E)$
 $\Rightarrow \text{int} * (E \text{ Op } E)$
 $\Rightarrow \text{int} * (\text{int} \text{ Op } E)$
 $\Rightarrow \text{int} * (\text{int} \text{ Op } \text{int})$
 $\Rightarrow \text{int} * (\text{int} + \text{int})$

Arithmetic Expressions

- Suppose we want to describe all legal arithmetic expressions using addition, subtraction, multiplication, and division.
- Here is one possible CFG:

$E \rightarrow \text{int}$
$E \rightarrow E \text{ Op } E$
$E \rightarrow (E)$
$\text{Op} \rightarrow +$
$\text{Op} \rightarrow -$
$\text{Op} \rightarrow *$
$\text{Op} \rightarrow /$

E
 $\Rightarrow E \text{ Op } E$
 $\Rightarrow E \text{ Op } \text{int}$
 $\Rightarrow \text{int } \text{Op } \text{int}$
 $\Rightarrow \text{int } / \text{int}$

Context-Free Grammars

- Formally, a context-free grammar is a collection of four items:
 - A set of **nonterminal symbols** (also called **variables**),
 - A set of **terminal symbols** (the **alphabet** of the CFG)
 - A set of **production rules** saying how each nonterminal can be replaced by a string of terminals and nonterminals, and
 - A **start symbol** (which must be a nonterminal) that begins the derivation.

E → **int**

E → **E Op E**

E → **(E)**

Op → **+**

Op → **-**

Op → *****

Op → **/**

Some CFG Notation

- In today's slides, capital letters in **Bold Red Uppercase** will represent nonterminals.
 - e.g. **A, B, C, D**
- Lowercase letters in **blue monospace** will represent terminals.
 - e.g. **t, u, v, w**
- Lowercase Greek letters in *gray italics* will represent arbitrary strings of terminals and nonterminals.
 - e.g. *α, γ, ω*
- You don't need to use these conventions on your own; just make sure whatever you do is readable. ☺

A Notational Shorthand

E → int

E → **E Op E**

E → (**E**)

Op → +

Op → -

Op → *

Op → /

A Notational Shorthand

E → **int** | **E Op E** | **(E)**

Op → **+** | **-** | ***** | **/**

Derivations

$E \rightarrow E \text{ Op } E \mid \text{int} \mid (E)$
$\text{Op} \rightarrow + \mid * \mid - \mid /$

E
 $\Rightarrow E \text{ Op } E$
 $\Rightarrow E \text{ Op } (E)$
 $\Rightarrow E \text{ Op } (E \text{ Op } E)$
 $\Rightarrow E * (E \text{ Op } E)$
 $\Rightarrow \text{int} * (E \text{ Op } E)$
 $\Rightarrow \text{int} * (\text{int} \text{ Op } E)$
 $\Rightarrow \text{int} * (\text{int} \text{ Op } \text{int})$
 $\Rightarrow \text{int} * (\text{int} + \text{int})$

- A sequence of steps where nonterminals are replaced by the right-hand side of a production is called a *derivation*.
- If string α derives string ω , we write $\alpha \Rightarrow^* \omega$.
- In the example on the left, we see $E \Rightarrow^* \text{int} * (\text{int} + \text{int})$.

The Language of a Grammar

- If G is a CFG with alphabet Σ and start symbol \mathbf{S} , then the *language of G* is the set

$$\mathcal{L}(G) = \{ \omega \in \Sigma^* \mid \mathbf{S} \Rightarrow^* \omega \}$$

- That is, $\mathcal{L}(G)$ is the set of strings derivable from the start symbol.
- Note: ω must be in Σ^* , the set of strings made from terminals. Strings involving nonterminals aren't in the language.

Context-Free Languages

- A language L is called a ***context-free language*** (or CFL) if there is a CFG G such that $L = \mathcal{L}(G)$.
- Questions:
 - What languages are context-free?
 - How are context-free and regular languages related?

From Regexes to CFGs

- CFGs consist purely of production rules of the form $A \rightarrow \omega$. They do not have the regular expression operators * or U.
- However, we can convert regular expressions to CFGs as follows:

$$S \rightarrow a^*b$$

From Regexes to CFGs

- CFGs consist purely of production rules of the form $A \rightarrow \omega$. They do not have the regular expression operators * or U.
- However, we can convert regular expressions to CFGs as follows:

$$S \rightarrow Ab$$

$$A \rightarrow Aa \mid \epsilon$$

From Regexes to CFGs

- CFGs consist purely of production rules of the form $A \rightarrow \omega$. They do not have the regular expression operators * or U.
- However, we can convert regular expressions to CFGs as follows:

$$S \rightarrow a (b \cup c^*)$$

From Regexes to CFGs

- CFGs consist purely of production rules of the form $A \rightarrow \omega$. They do not have the regular expression operators $*$ or \cup .
- However, we can convert regular expressions to CFGs as follows:

$$S \rightarrow aX$$

$$X \rightarrow b \mid C$$

$$C \rightarrow Cc \mid \epsilon$$

Regular Languages and CFLs

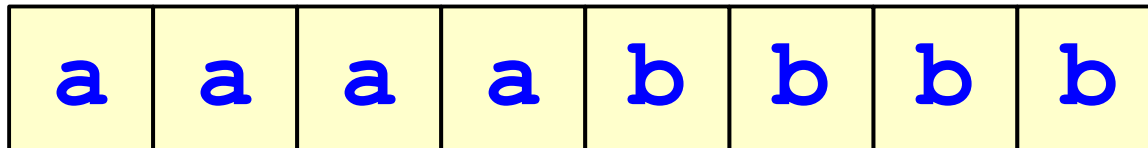
- ***Theorem:*** Every regular language is context-free.
- ***Proof Idea:*** Use the construction from the previous slides to convert a regular expression for L into a CFG for L . ■
- ***Problem Set 8 Exercise:*** Instead, show how to convert a DFA/NFA into a CFG.

The Language of a Grammar

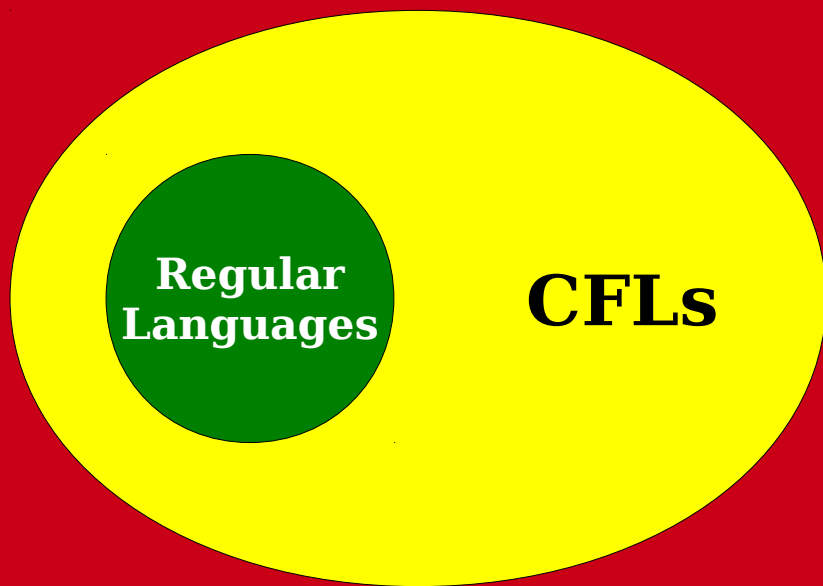
- Consider the following CFG G :

$$S \rightarrow aSb \mid \epsilon$$

- What strings can this generate?



$$\mathcal{L}(G) = \{ a^n b^n \mid n \in \mathbb{N} \}$$

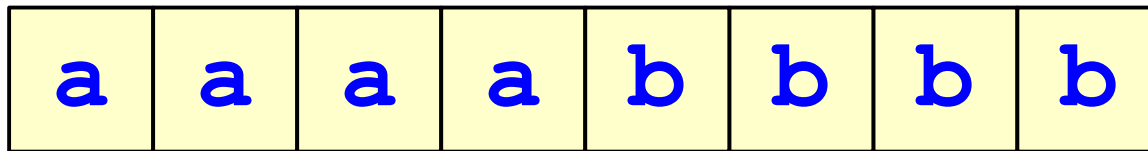


All Languages

Why the Extra Power?

- Why do CFGs have more power than regular expressions?
- ***Intuition:*** Derivations of strings have unbounded “memory.”

$$S \rightarrow aSb \mid \epsilon$$



Time-Out for Announcements!

Problem Sets

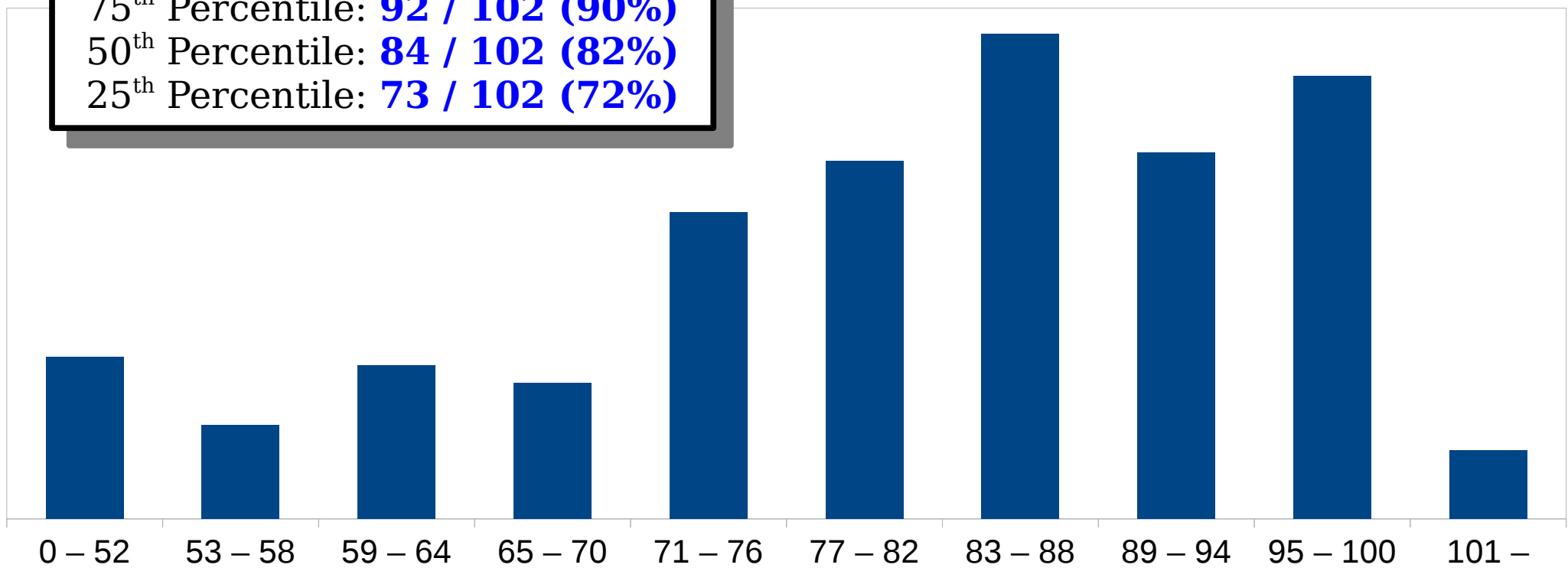
- Problem Set Six was due at 2:30PM today.
 - ***Be careful about using late days here.*** The midterm is on Monday.
- Problem Set Seven goes out today. It's due next Friday at 2:30PM.
 - Play around with regular expressions, the power of regular languages, and the limits of regularity!

Problem Set 5 Grades

75th Percentile: **92 / 102 (90%)**

50th Percentile: **84 / 102 (82%)**

25th Percentile: **73 / 102 (72%)**

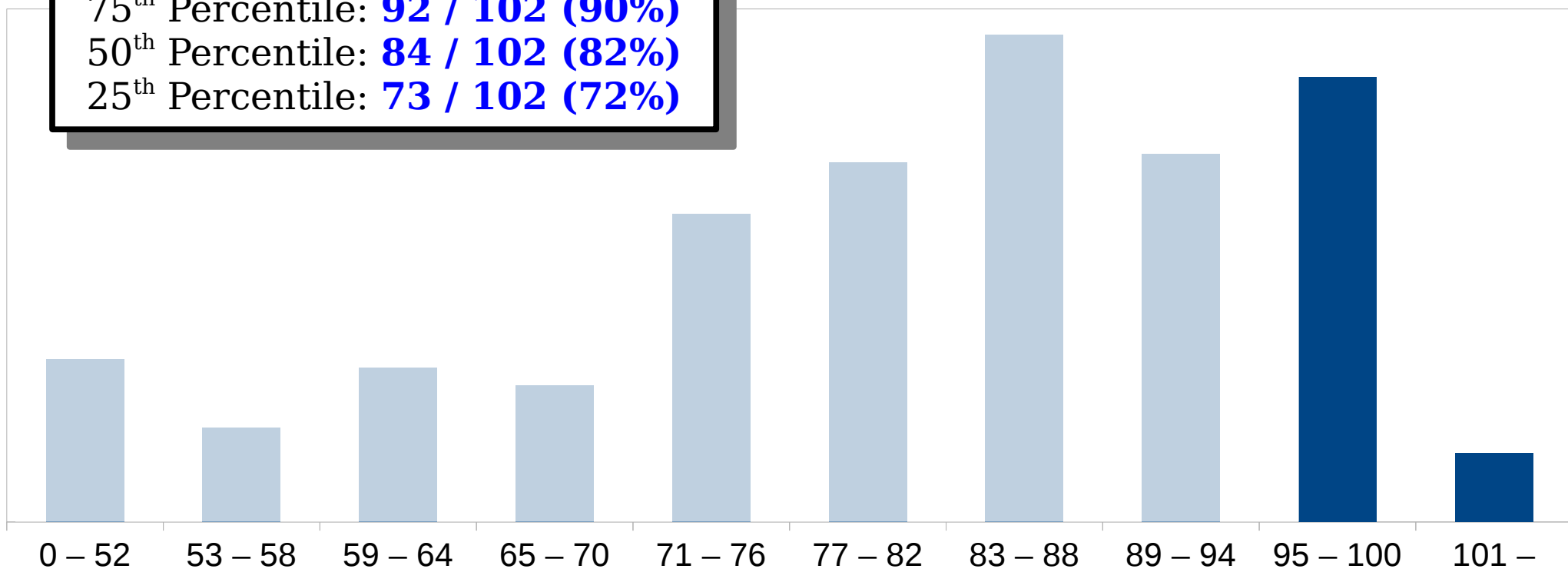


Problem Set 5 Grades

75th Percentile: **92 / 102 (90%)**

50th Percentile: **84 / 102 (82%)**

25th Percentile: **73 / 102 (72%)**



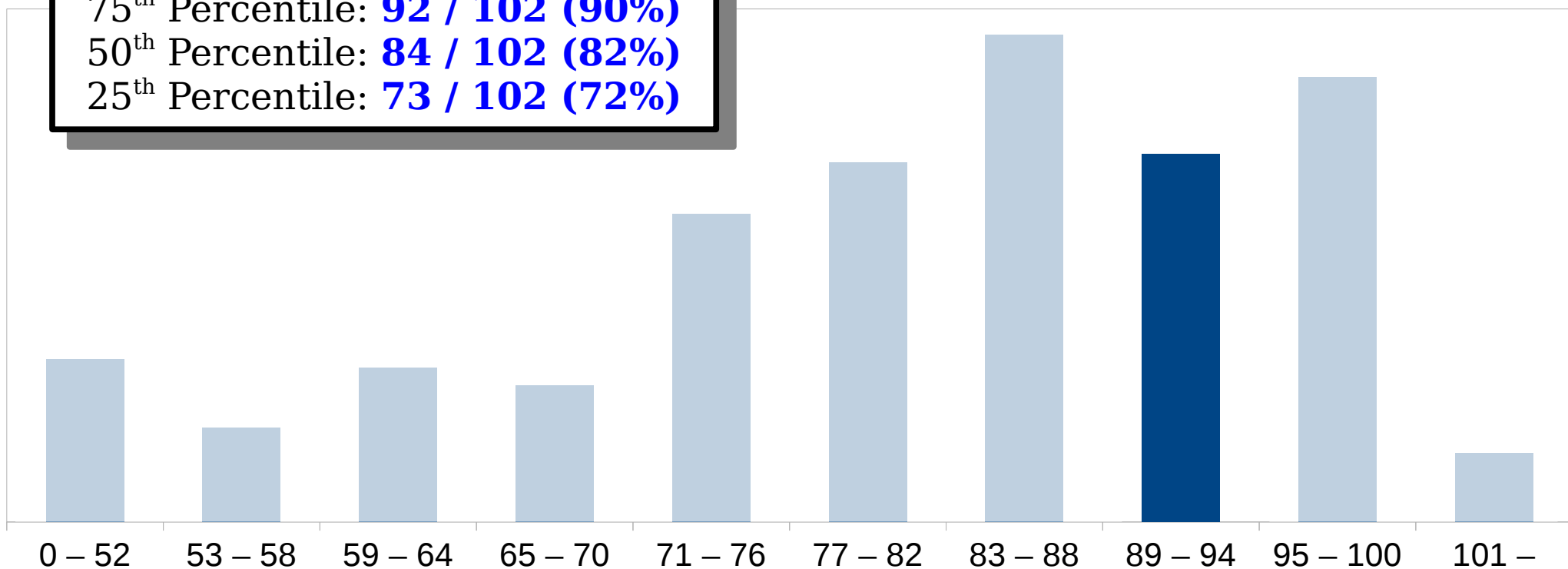
Awesome job! Take a look at your feedback to see how to patch up those last few areas.

Problem Set 5 Grades

75th Percentile: **92 / 102 (90%)**

50th Percentile: **84 / 102 (82%)**

25th Percentile: **73 / 102 (72%)**



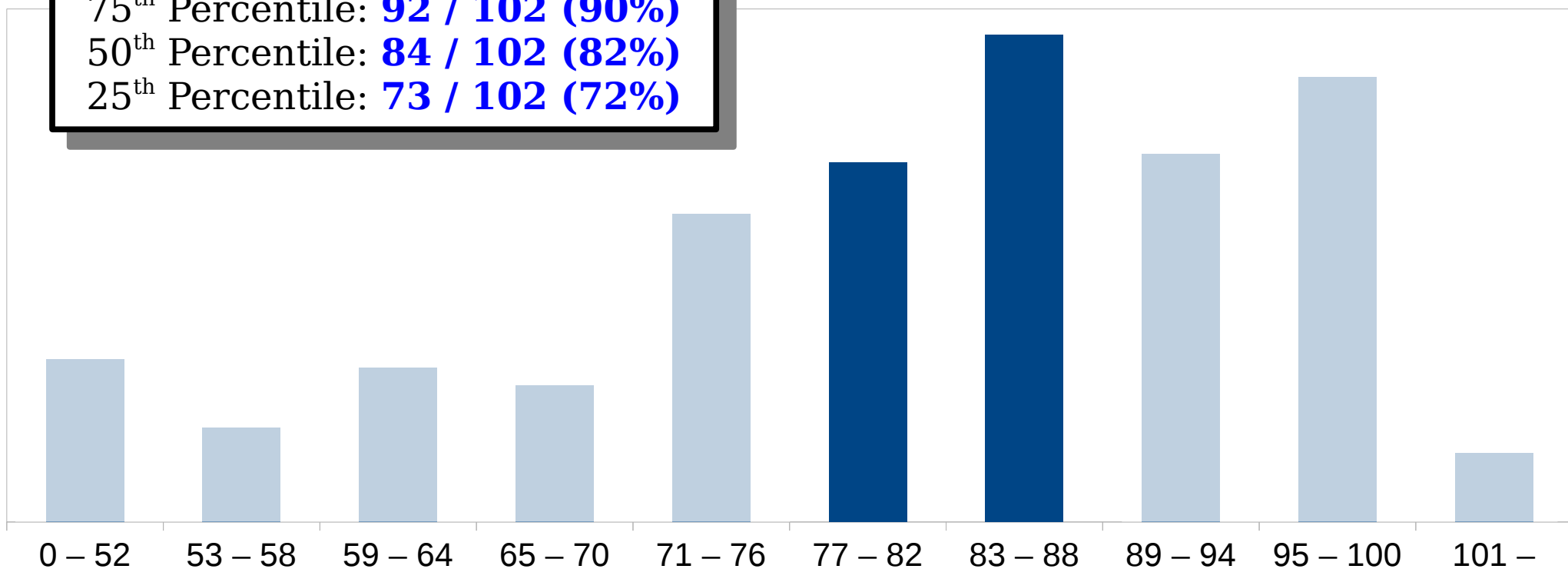
Well done! Review your feedback to see if there are any repeated issues you're running into.

Problem Set 5 Grades

75th Percentile: **92 / 102 (90%)**

50th Percentile: **84 / 102 (82%)**

25th Percentile: **73 / 102 (72%)**



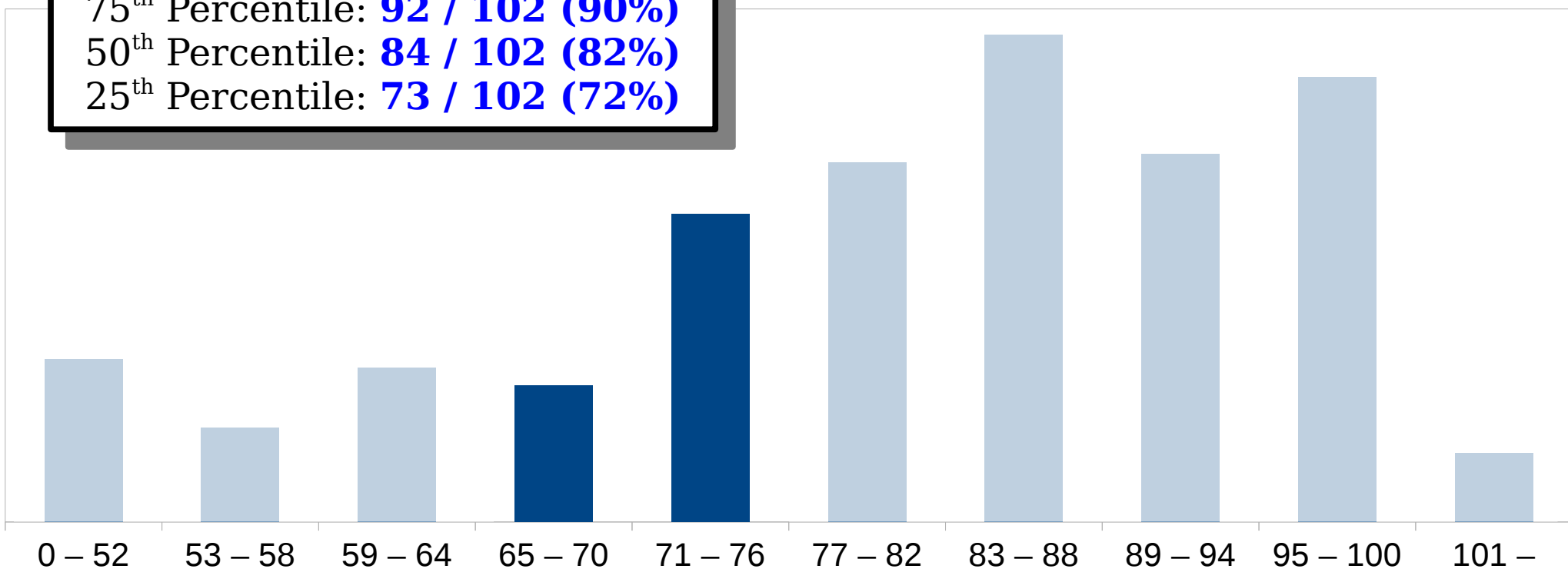
A solid performance! Take a look at your feedback, redo any problems you had trouble with, and check the solutions set if you got really tripped up here.

Problem Set 5 Grades

75th Percentile: **92 / 102 (90%)**

50th Percentile: **84 / 102 (82%)**

25th Percentile: **73 / 102 (72%)**



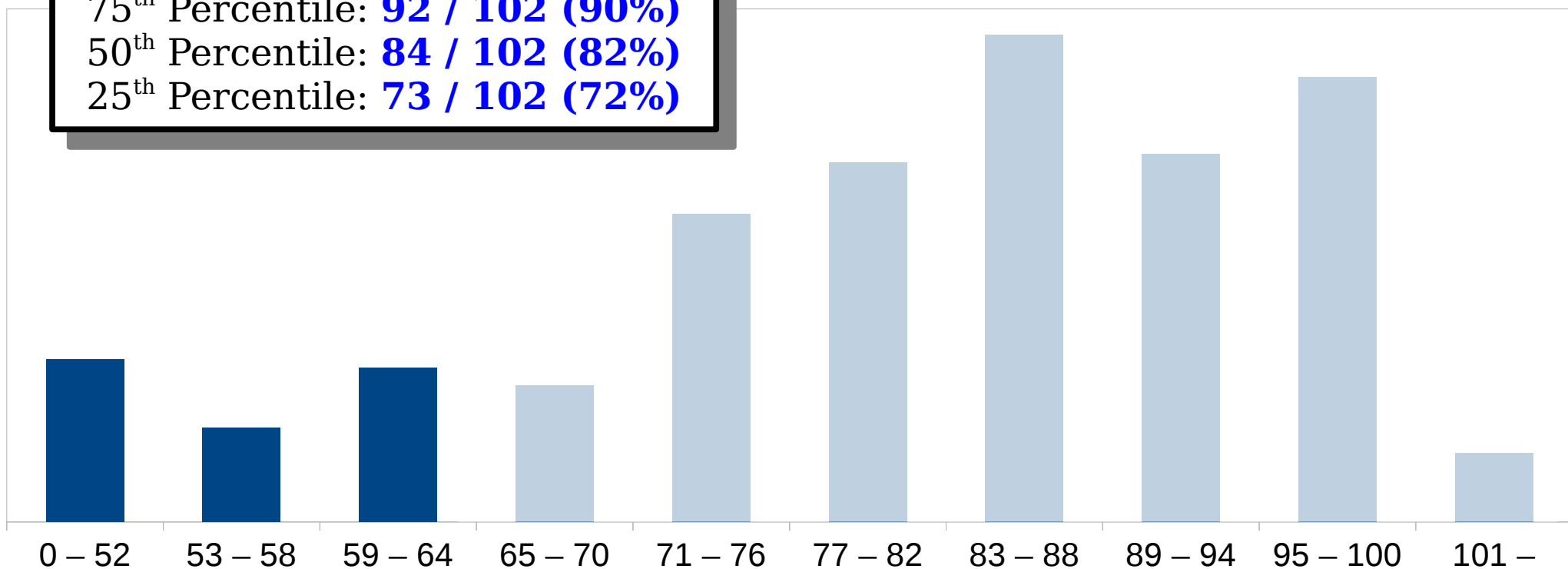
You're on the right track, but it looks like there are some concepts that haven't yet clicked. Look at your grader's feedback, redo problems that you had trouble with, and stop by office hours / ask questions on Piazza on topics you're not fully solidified on.

Problem Set 5 Grades

75th Percentile: **92 / 102 (90%)**

50th Percentile: **84 / 102 (82%)**

25th Percentile: **73 / 102 (72%)**



Looks like induction is giving you some trouble. Stop by office hours to get some input and feedback on how to improve, and work through some of the extra practice problems on induction to solidify things in advance of the exam.

Midterm Exam Logistics

- The next midterm is **Monday, November 13th** from **7:00PM - 10:00PM**. Locations are divvied up by last (family) name:
 - Abb - Hal: Go to **Hewlett 201**.
 - Han - Zwa: Go to **Hewlett 200**.
- The exam focuses on Lecture 06 – 13 (binary relations through induction) and PS3 – PS5. Finite automata onward is *not* tested.
 - Topics from earlier in the quarter (proofwriting, first-order logic, set theory, etc.) are also fair game, but that's primarily because the later material builds on this earlier material.
- The exam is closed-book, closed-computer, and limited-note. You can bring a double-sided, 8.5" × 11" sheet of notes with you to the exam, decorated however you'd like.

Practice Materials

- We have a ton of practice materials up on the course website:
 - EPP2: a collection of 21 practice problems on a variety of topics.
 - Four practice midterms, with solutions.
 - Only have time for one? Do **Practice Second Midterm 4**. That's the one we did at the practice exam on Wednesday.
- And please feel free to ask questions on Piazza over the next couple of days – we're happy to help out!

Three Questions

- What is something you know now that, at the start of the quarter, you knew you didn't know?
- What is something you know now that, at the start of the quarter, you *didn't* know that you didn't know?
- What is something you *don't* know that, at the start of the quarter, you *didn't* know that you didn't know?

Your Questions

“Do you have any suggestions for interesting topics in CS to wikipedia? I don't know what I don't know”

Wikipedia is not necessarily the best resource to pick up new CS topics, since it's often written assuming you already know what's out there. I'd look at places like Coursera, Udacity, edX, and Khan Academy as launching points for learning new topics. Having a quick intro to the subject makes it a lot easier to digest new topics!

“What should we do if we're not sure what grader comments want us to change about our work?”

Feel free to stop by office hours or to ask on Piazza if you need a clarification. We're happy to help out!

“Why is it still so hard to find a job/internship when I hear from many CS professors that the demand for software engineer jobs/positions is increasing?”

To really answer this question, I'd need to hear about where you are (freshman? grad student?) and what your background experience is. You have an open offer to stop by my office hours to chat about this or to email me to set up a time to meet to talk about this!

“How do you define success? How do I succeed?”

I like to think of “success” as a limit rather than achievable goal. No matter what you do, there’s always more you can be doing and there’s always room for improvement. Try to do the best you can about the things that matter to you, and make sure that you have a broad rather than narrow focus of what you think is important.

“why”

— \ — (ツ) — / —

Back to CS103!

Designing CFGs

- Like designing DFAs, NFAs, and regular expressions, designing CFGs is a craft.
- When thinking about CFGs:
 - ***Think recursively:*** Build up bigger structures from smaller ones.
 - ***Have a construction plan:*** Know in what order you will build up the string.
 - ***Store information in nonterminals:*** Have each nonterminal correspond to some useful piece of information.

Designing CFGs

- Let $\Sigma = \{a, b\}$ and let $L = \{w \in \Sigma^* \mid w \text{ is a palindrome}\}$
- We can design a CFG for L by thinking inductively:
 - Base case: ϵ , a , and b are palindromes.
 - If w is a palindrome, then awa and bwb are palindromes.
 - No other strings are palindromes.

S \rightarrow ϵ | **a** | **b** | **aSa** | **bSb**

Designing CFGs

- Let $\Sigma = \{ (,) \}$ and let $L = \{ w \in \Sigma^* \mid w \text{ is a string of balanced parentheses} \}$
- Some sample strings in L :

$((()))$

$(())()$

$((()))((()))$

$((((()))(()))$

ϵ

$()()$

Designing CFGs

- Let $\Sigma = \{ (,) \}$ and let $L = \{ w \in \Sigma^* \mid w \text{ is a string of balanced parentheses} \}$
- Let's think about this recursively.
 - Base case: the empty string is a string of balanced parentheses.
 - Recursive step: Look at the closing parenthesis that matches the first open parenthesis.

The diagram shows a string of blue parentheses: $((()((())))((()((())))$. A vertical dashed red line is positioned between the closing parenthesis of the first pair and the opening parenthesis of the second pair, illustrating the recursive step of identifying the first matching pair.

Designing CFGs

- Let $\Sigma = \{ (,) \}$ and let $L = \{ w \in \Sigma^* \mid w \text{ is a string of balanced parentheses} \}$
- Let's think about this recursively.
 - Base case: the empty string is a string of balanced parentheses.
 - Recursive step: Look at the closing parenthesis that matches the first open parenthesis. Removing the first parenthesis and the matching parenthesis forms two new strings of balanced parentheses.

$$S \rightarrow (S)S \mid \epsilon$$

Designing CFGs

- Let $\Sigma = \{a, b\}$ and let $L = \{w \in \Sigma^* \mid w \text{ has the same number of } a\text{'s and } b\text{'s}\}$
- Is this a CFG for L ?

$$S \rightarrow aSb \mid bSa \mid \epsilon$$

- Can you derive the string **abba**?

Designing CFGs: A Caveat

- When designing a CFG for a language, make sure that it
 - generates all the strings in the language and
 - never generates a string outside the language.
- The first of these can be tricky - make sure to test your grammars!
- You'll design your own CFG for this language on Problem Set 8.

CFG Caveats II

- Is the following grammar a CFG for the language $\{ a^n b^n \mid n \in \mathbb{N} \}$?

$$S \rightarrow aSb$$

- What strings in $\{a, b\}^*$ can you derive?
 - Answer: ***None!***
- What is the language of the grammar?
 - Answer: \emptyset
- When designing CFGs, make sure your recursion actually terminates!

CFG Caveats III

- When designing CFGs, remember that each nonterminal can be expanded out independently of the others.
- Let $\Sigma = \{a, \underline{a}\}$ and let $L = \{a^n \underline{a} a^n \mid n \in \mathbb{N}\}$.
- Is the following a CFG for L ?

$$S \rightarrow X \underline{a} X$$

$$X \rightarrow aX \mid \varepsilon$$

$$\begin{aligned} S &\Rightarrow X \underline{a} X \\ &\Rightarrow aX \underline{a} X \\ &\Rightarrow aaX \underline{a} X \\ &\Rightarrow aa \underline{a} X \\ &\Rightarrow aa \underline{a} aX \\ &\Rightarrow aa \underline{a} a \end{aligned}$$

Finding a Build Order

- Let $\Sigma = \{a, \stackrel{?}{=}\}$ and let $L = \{a^n \stackrel{?}{=} a^n \mid n \in \mathbb{N}\}$.
- To build a CFG for L , we need to be more clever with how we construct the string.
 - If we build the strings of a 's independently of one another, then we can't enforce that they have the same length.
 - **Idea:** Build both strings of a 's at the same time.
- Here's one possible grammar based on that idea:

$$S \rightarrow \stackrel{?}{=} \mid aSa$$

S

$$\Rightarrow aSa$$

$$\Rightarrow aaSaa$$

$$\Rightarrow aaaSaaa$$

$$\Rightarrow aaa \stackrel{?}{=} aaa$$

Function Prototypes

- Let $\Sigma = \{\text{void, int, double, name, (,) , , , ;}\}$.
- Let's write a CFG for C-style function prototypes!
- Examples:
 - `void name(int name, double name);`
 - `int name();`
 - `int name(double name);`
 - `int name(int, int name, int);`
 - `void name(void);`

Function Prototypes

- Here's one possible grammar:
 - **S** → **Ret** **name** (**Args**) ;
 - **Ret** → **Type** | **void**
 - **Type** → **int** | **double**
 - **Args** → ϵ | **void** | **ArgList**
 - **ArgList** → **OneArg** | **ArgList**, **OneArg**
 - **OneArg** → **Type** | **Type** **name**
- Fun question to think about: what changes would you need to make to support pointer types?

Summary of CFG Design Tips

- Look for recursive structures where they exist: they can help guide you toward a solution.
- Keep the build order in mind – often, you'll build two totally different parts of the string concurrently.
 - Usually, those parts are built in opposite directions: one's built left-to-right, the other right-to-left.
- Use different nonterminals to represent different structures.

Applications of Context-Free Grammars

CFGs for Programming Languages

BLOCK → **STMT**
 | **{ STMTS }**

STMTS → ϵ
 | **STMT STMTS**

STMT → **EXPR;**
 | **if (EXPR) BLOCK**
 | **while (EXPR) BLOCK**
 | **do BLOCK while (EXPR);**
 | **BLOCK**
 | ...

EXPR → **identifier**
 | **constant**
 | **EXPR + EXPR**
 | **EXPR - EXPR**
 | **EXPR * EXPR**
 | ...

Grammars in Compilers

- One of the key steps in a compiler is figuring out what a program “means.”
- This is usually done by defining a grammar showing the high-level structure of a programming language.
- There are certain classes of grammars (LL(1) grammars, LR(1) grammars, LALR(1) grammars, etc.) for which it's easy to figure out how a particular string was derived.
- Tools like yacc or bison automatically generate parsers from these grammars.
- Curious to learn more? Take CS143!

Natural Language Processing

- By building context-free grammars for actual languages and applying statistical inference, it's possible for a computer to recover the likely meaning of a sentence.
 - In fact, CFGs were first called ***phrase-structure grammars*** and were introduced by Noam Chomsky in his seminal work *Syntactic Structures*.
 - They were then adapted for use in the context of programming languages, where they were called ***Backus-Naur forms***.
- Stanford's **CoreNLP project** is one place to look for an example of this.
- Want to learn more? Take CS124 or CS224N!

Next Time

- ***No class on Monday - you have a midterm!***
- ***Then, when we get back...***
 - ***Turing Machines***
 - What does a computer with unbounded memory look like?
 - How do you program them?