The background features a dark purple gradient with a subtle geometric pattern. On the left side, there is a dense cluster of white circular nodes connected by thin blue lines, forming a network graph. To the right of this cluster, several light blue triangular outlines are scattered across the slide.

Stack Overflow NLP Predictions

Vanessa Alvarado

Background	01
Problem Statement	02
Data Overview	03
TABLE OF CONTENTS	
Modeling	04
Conclusion	05
Next Steps	06

01

Background



Stack Overflow

Stack Overflow is an online community for developers to learn and share their programming knowledge.

A complex network graph is visible in the background, composed of numerous small white dots connected by thin white lines, creating a web-like structure.

02

Problem Statement

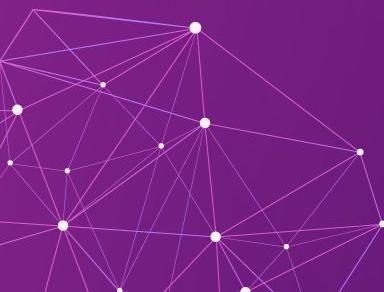
In 2019...

2.8

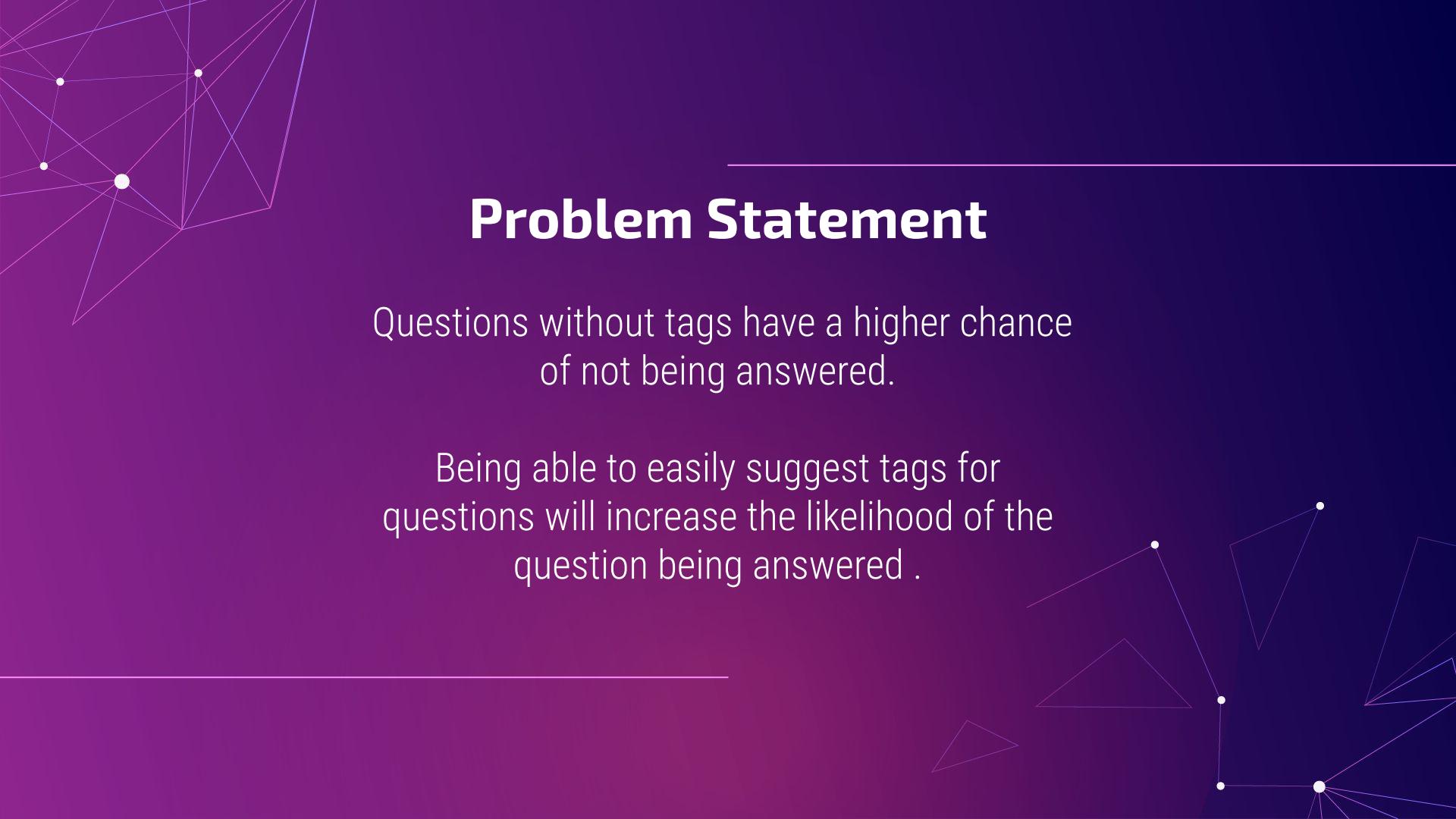
Million Answers

2.6

Million New Questions



* Stack Overflow



Problem Statement

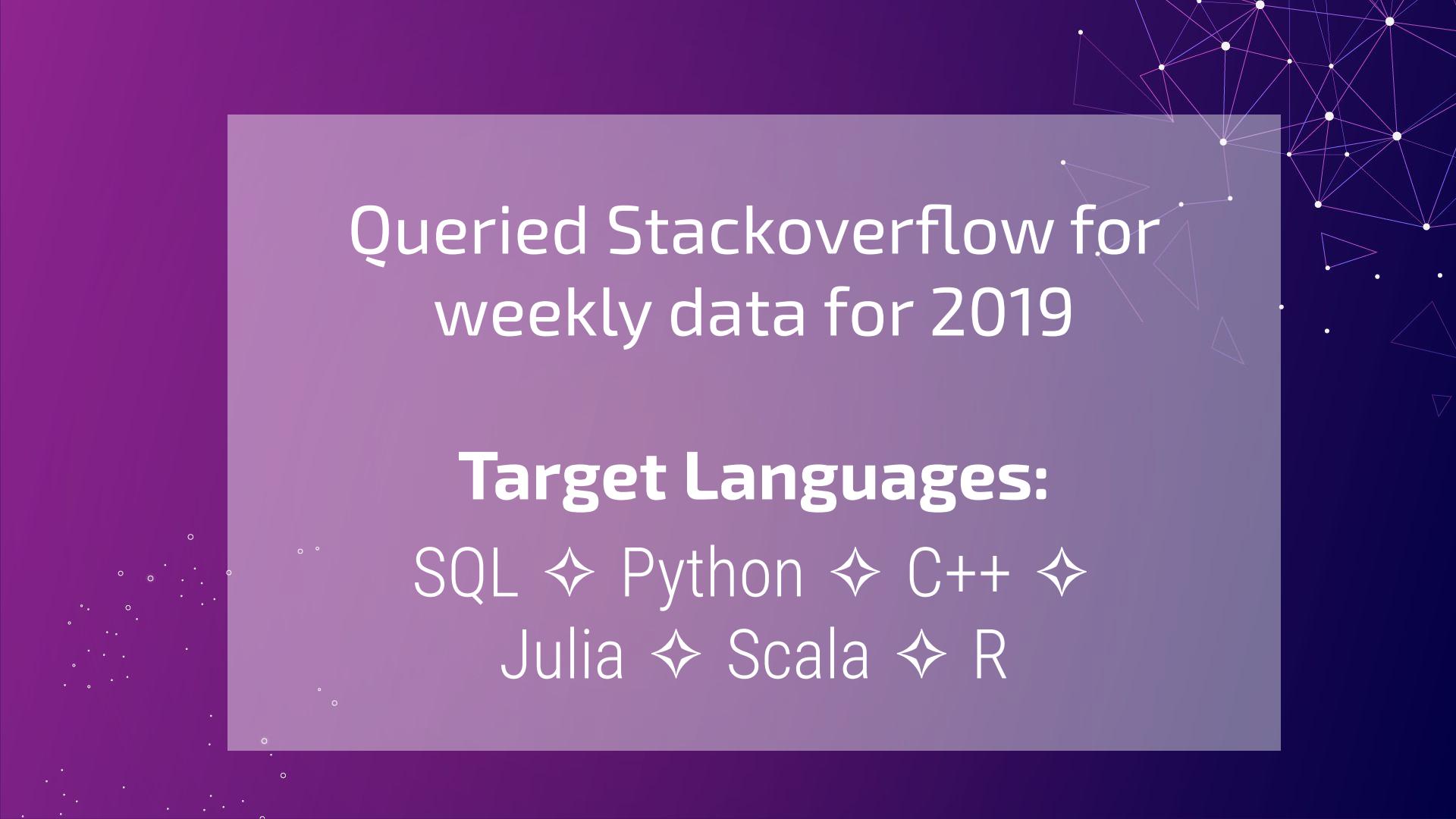
Questions without tags have a higher chance of not being answered.

Being able to easily suggest tags for questions will increase the likelihood of the question being answered .



03

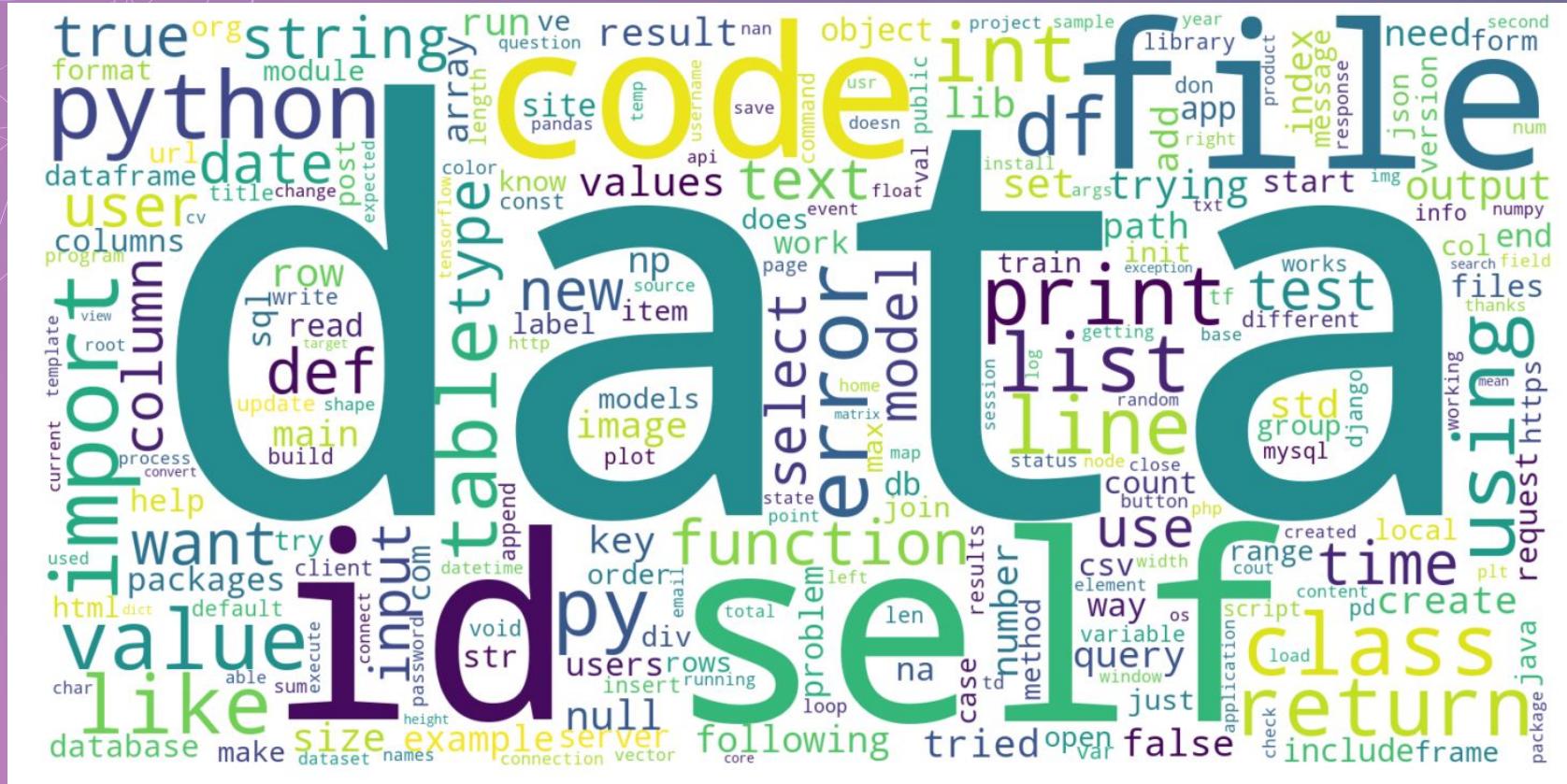
Data Overview

A dark purple background featuring a faint, abstract network graph composed of small white dots and thin white lines.

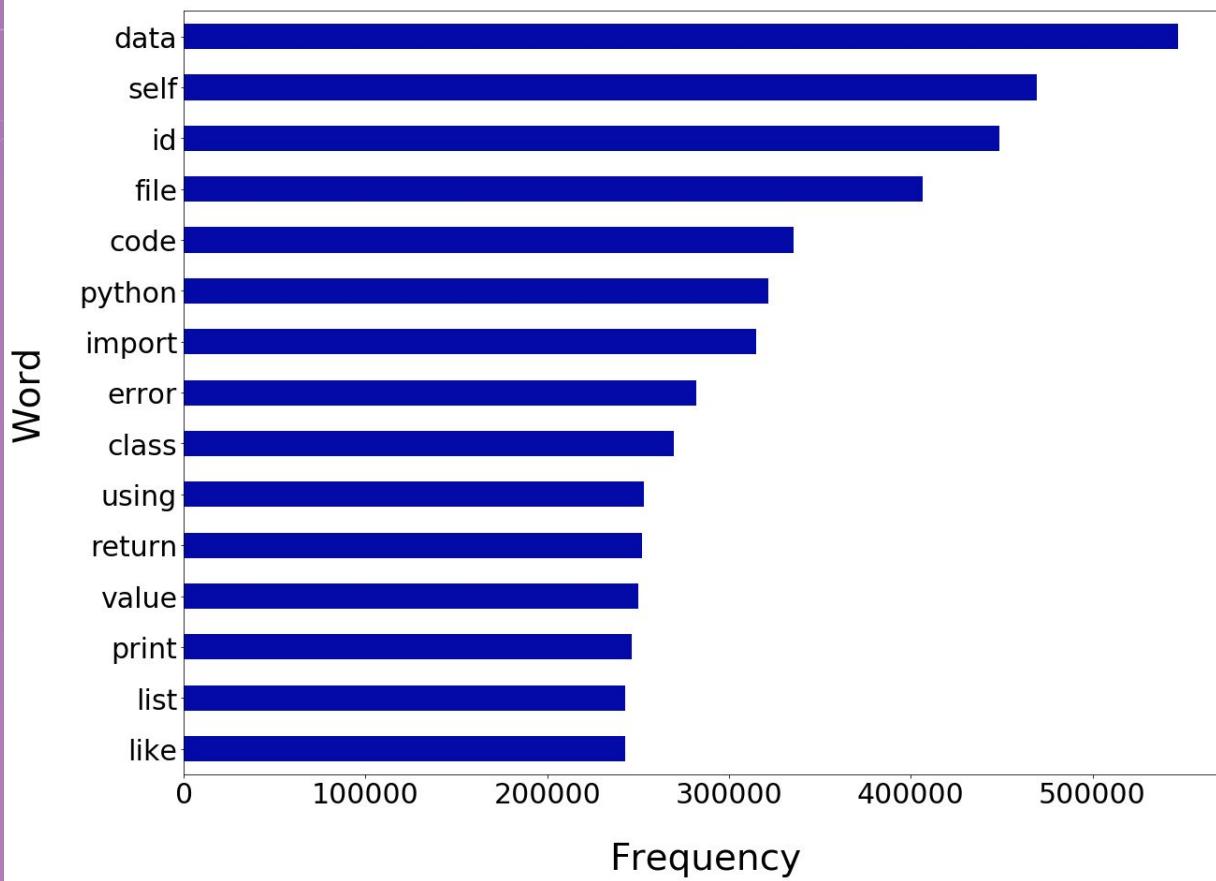
Queried Stackoverflow for weekly data for 2019

Target Languages:

SQL ✦ Python ✦ C++ ✦
Julia ✦ Scala ✦ R



Top 15 Words: Total



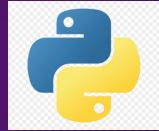
Language Specific Top Words

SQL



1. Id
2. Table
3. Select
4. Data
5. Query

Python



1. Self
2. File
3. Python
4. Import
5. Data

Julia



1. Julia
2. Jl
3. End
4. Function
5. Array

R



1. Data
2. Na
3. Df
4. Function
5. List

C++



1. Int
2. Std
3. Include
4. Return
5. Cout

Scala



1. Scala
2. Val
3. Spark
4. Org
5. String

Top Words Used Across Languages

Code

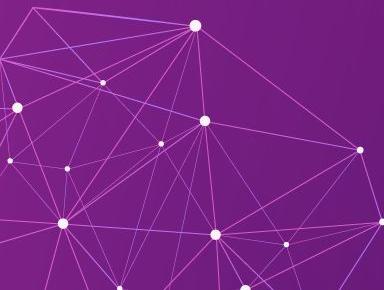
SQL, C++, Julia, R, Python

Error

SQL, C++, Julia, R, Scala

Data

SQL, Julia, R, Python



04

Modeling



Model Accuracy

Baseline: 55.23%

	Training	Testing	Diff.
Logistic + CountVect	89.96%	89.68%	0.28%
Logistic + TFIDF	91.43%	91.25%	0.18%
Decision Trees + TFIDF	82.99%	82.46%	0.53%



05

Conclusions



55% Python

21% SQL

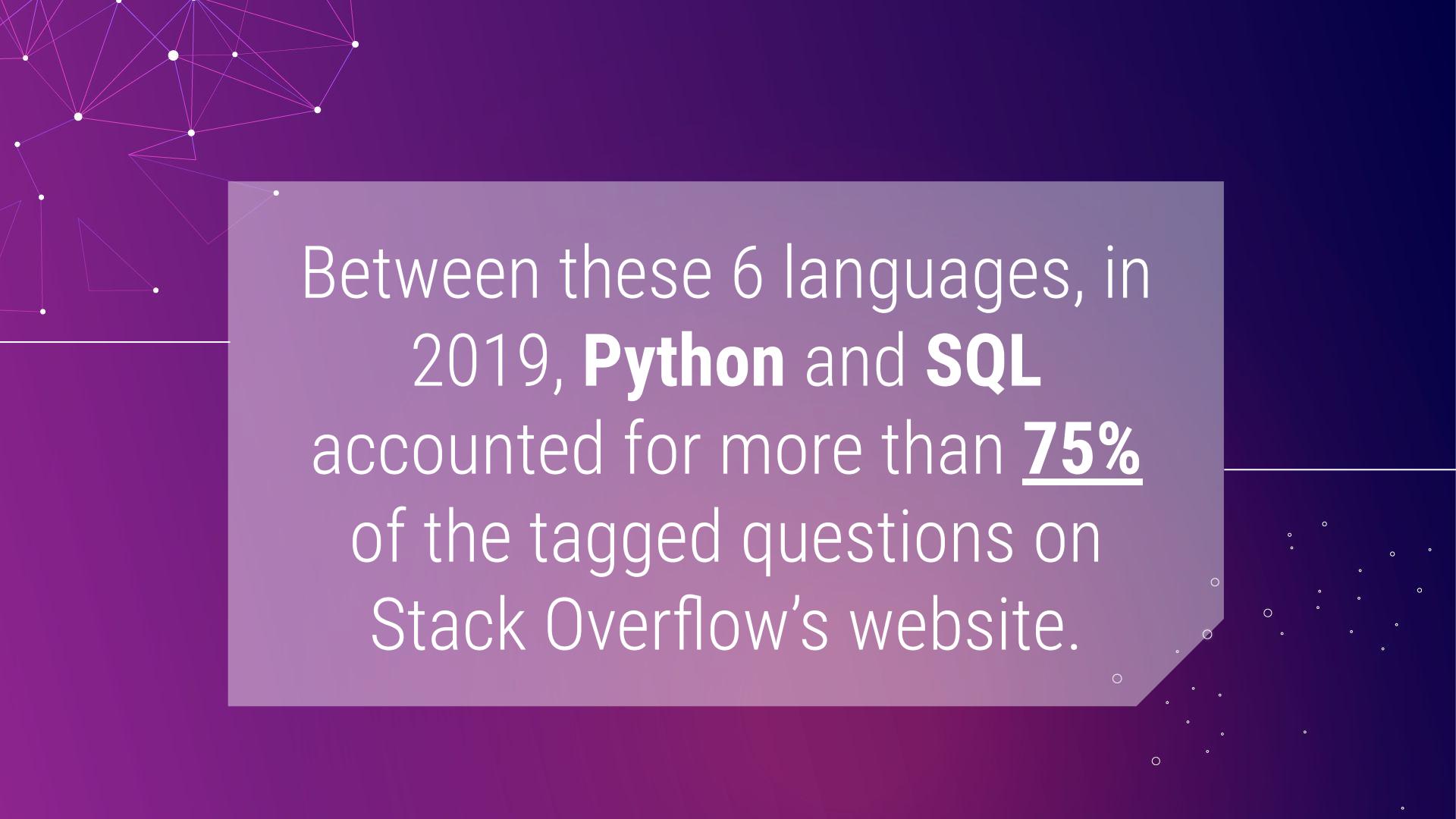
11% R

Question %

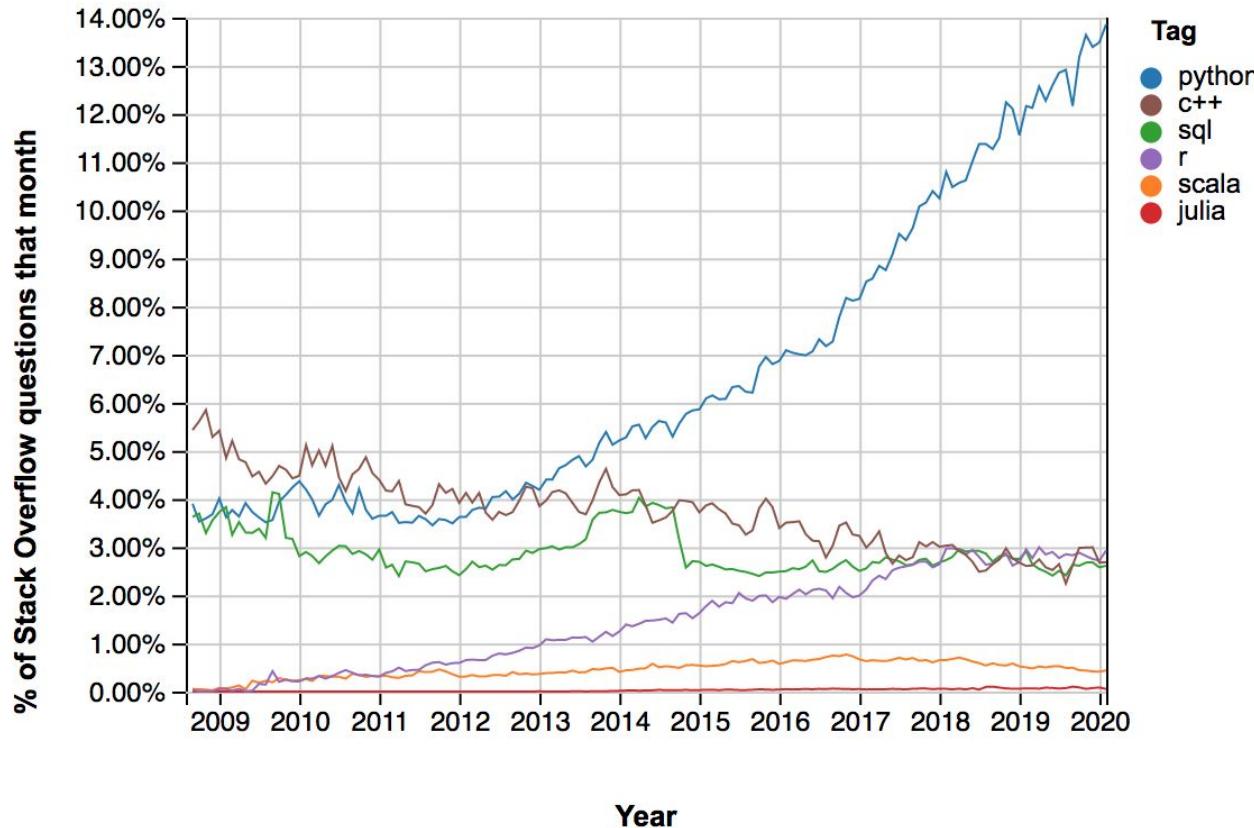
10% C++

1.9% Scala

0.3% Julia



Between these 6 languages, in
2019, **Python** and **SQL**
accounted for more than **75%**
of the tagged questions on
Stack Overflow's website.



06

Next Steps

Future Work

Data

Add more years to data used for modeling

More coding languages

Modeling

Hyperparameters

Random Forests

Word2Vect

Global Vectors (GloVe)

Resources

Cloud Computing to help fit and run more computationally heavy models.



*“The world is one big
data problem.”*

– Andrew McAfee



THANKS

Does anyone have
any questions?

CREDITS: This presentation template was created by [Slidesgo](#), including icons by [Flaticon](#), and infographics & images by [Freepik](#).

Please keep this slide for attribution.

RESOURCES

VECTORS

- Technology background with gradient colors
- Blue 5g concept background
- Abstract landing pages with technology devices

PHOTOS

- Woman using smartphone with hologram
- Happy businesswoman looking at camera with holding pencil and diary
- Portrait of smiling man holding digital tablet looking at camera
- Smiling bearded man holding disposable coffee cup while opening door
- Portrait of pretty woman holding laptop looking at camera
- Motherboard with optical fiber

