

TypyBench: Evaluating LLM Type Inference for Untyped Python Repositories

ICML 2025

Honghua Dong^{1,2*}, Jiacheng Yang^{1,2*}, Xun Deng^{1*}, Yuhe Jiang^{1,2*},
Gennady Pekhimenko^{1,2,3}, Fan Long¹, and Xujie Si^{1,2,3}

University of Toronto¹, Vector Institute² and CIFAR AI Chair³

*Equal Contribution

Introduction

- **Advantages of Type Annotations:**
 - Enhance code clarity
 - Prevent type-related errors
 - Enable autocompletion
 - Facilitate maintenance & refactoring

```
class Greeter:  
    def __init__(self, name: str) → None:  
        self.name = name  
    def greet(self) → str:  
        return f"Hello, {self.name}"  
    def get_name(self) → str:  
        return self.name
```

Introduction

- **Advantages of Type Annotations**
- **Type Inference**
 - Given code, infer types

Untyped Code

```
class Greeter:  
    def __init__(self, name):  
        self.name = name  
    def greet(self):  
        return f"Hello, {self.name}"  
    def get_name(self):  
        return self.name
```

Introduction

- Advantages of Type Annotations
- Type Inference
 - Given code, infer types

Type Annotated Code

```
class Greeter:  
    def __init__(self, name: str) → None:  
        self.name = name  
    def greet(self) → str:  
        return f"Hello, {self.name}"  
    def get_name(self) → str:  
        return self.name
```

Introduction

- Advantages of Type Annotations
- Type Inference
- Benchmarks

Existing Benchmarks

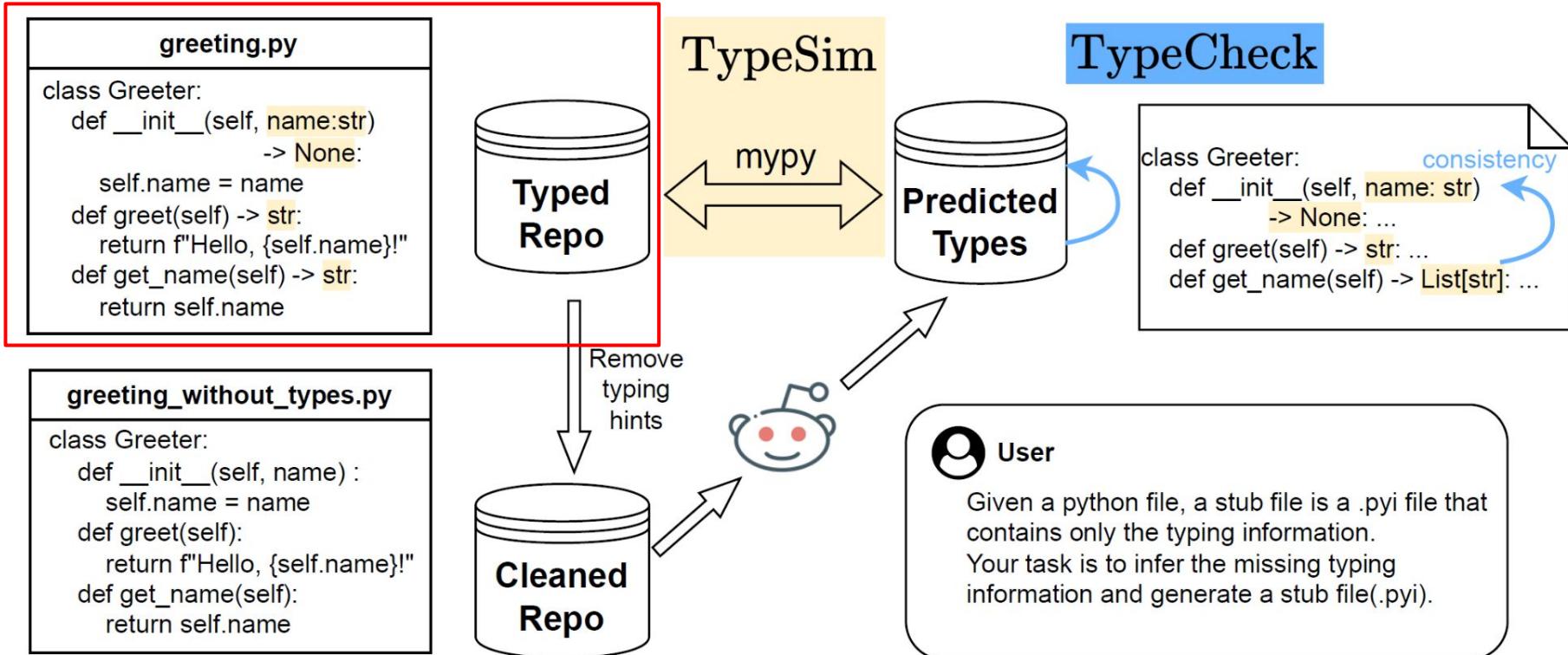
- ✗ Mostly rely on exact matching
- ✗ Less attention on type consistency
- ✗ No large-scale evaluation on LLMs

Introduction

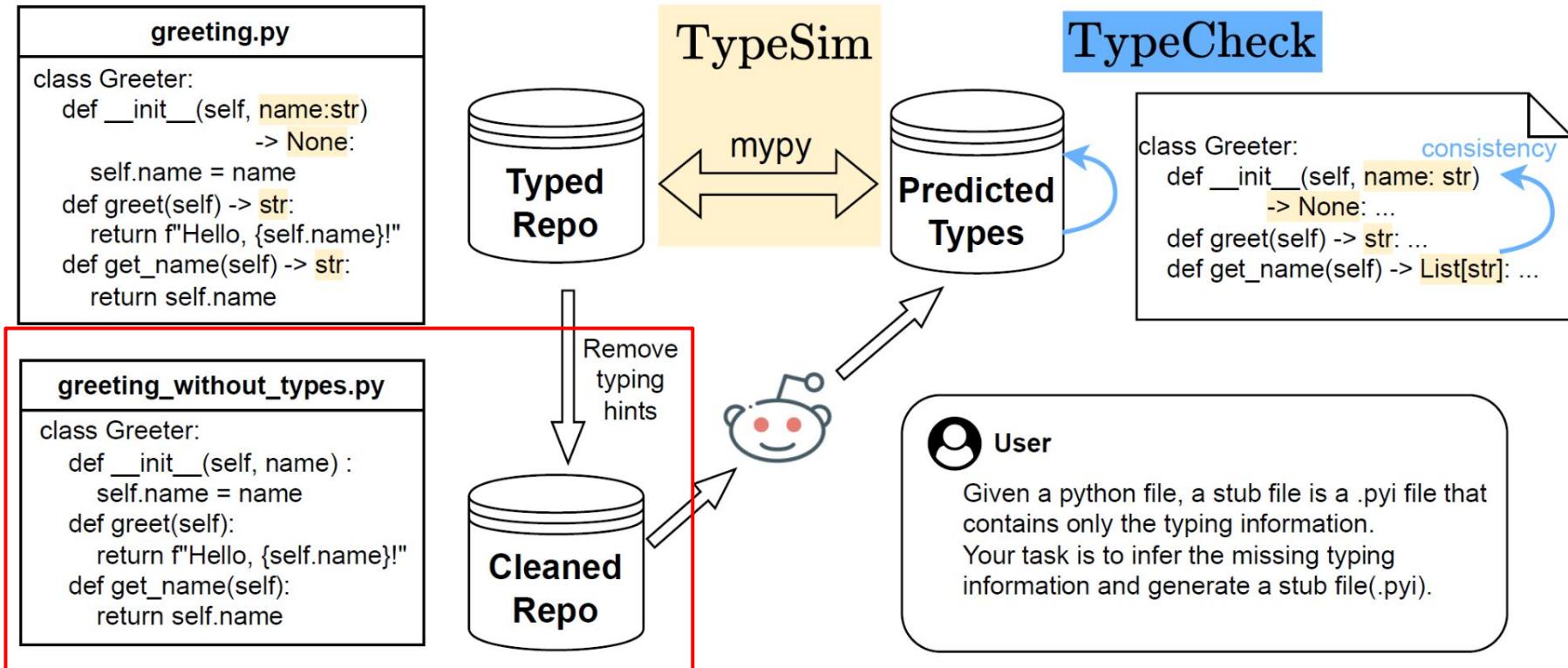
- Advantages of Type Annotations
- Type Inference
- Benchmarks

Existing Benchmarks	TypyBench
✗ Mostly rely on exact matching	✓ TypeSim Metric
✗ Less attention on type consistency	✓ TypeCheck Metric
✗ No large-scale evaluation on LLMs	✓ Holistic Evaluation of SOTA LLMs

TypyBench



TypyBench



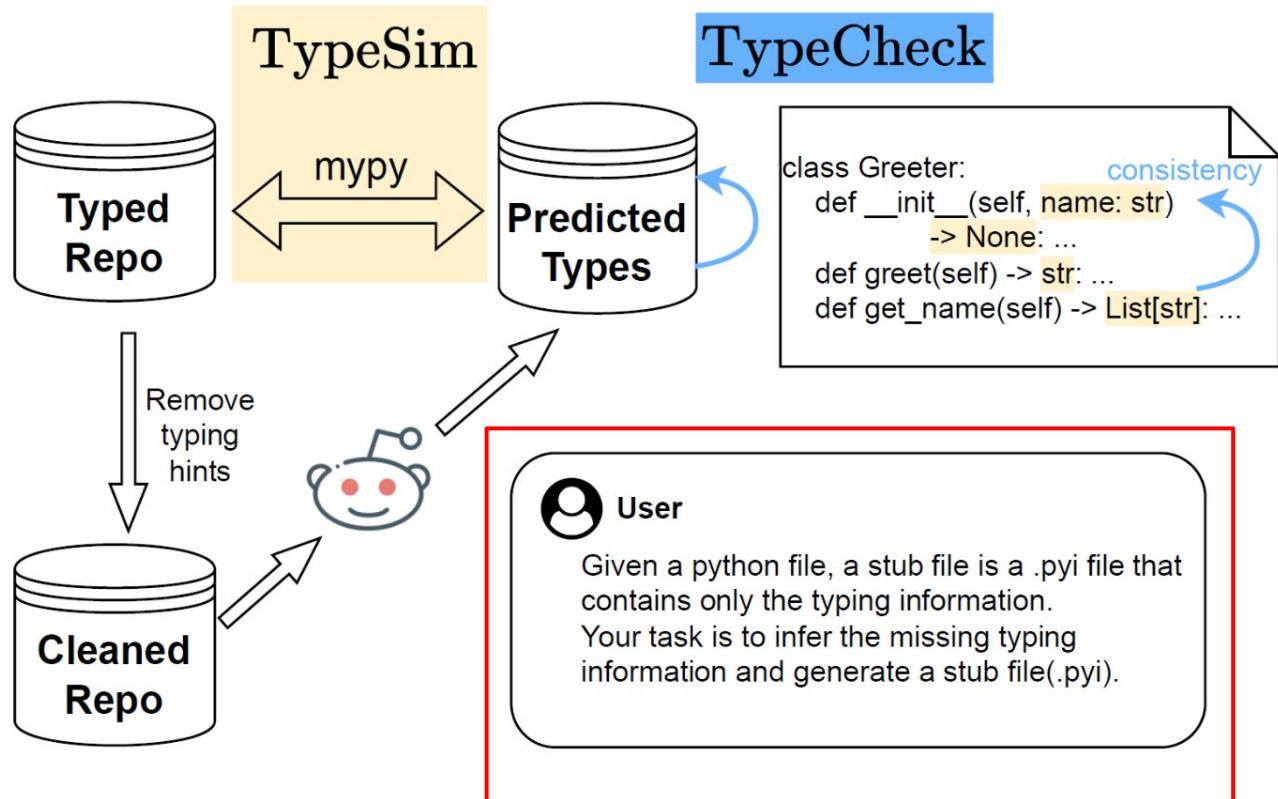
TypyBench

```
greeting.py
```

```
class Greeter:  
    def __init__(self, name:str)  
        -> None:  
        self.name = name  
    def greet(self) -> str:  
        return f"Hello, {self.name}!"  
    def get_name(self) -> str:  
        return self.name
```

```
greeting_without_types.py
```

```
class Greeter:  
    def __init__(self, name):  
        self.name = name  
    def greet(self):  
        return f"Hello, {self.name}!"  
    def get_name(self):  
        return self.name
```



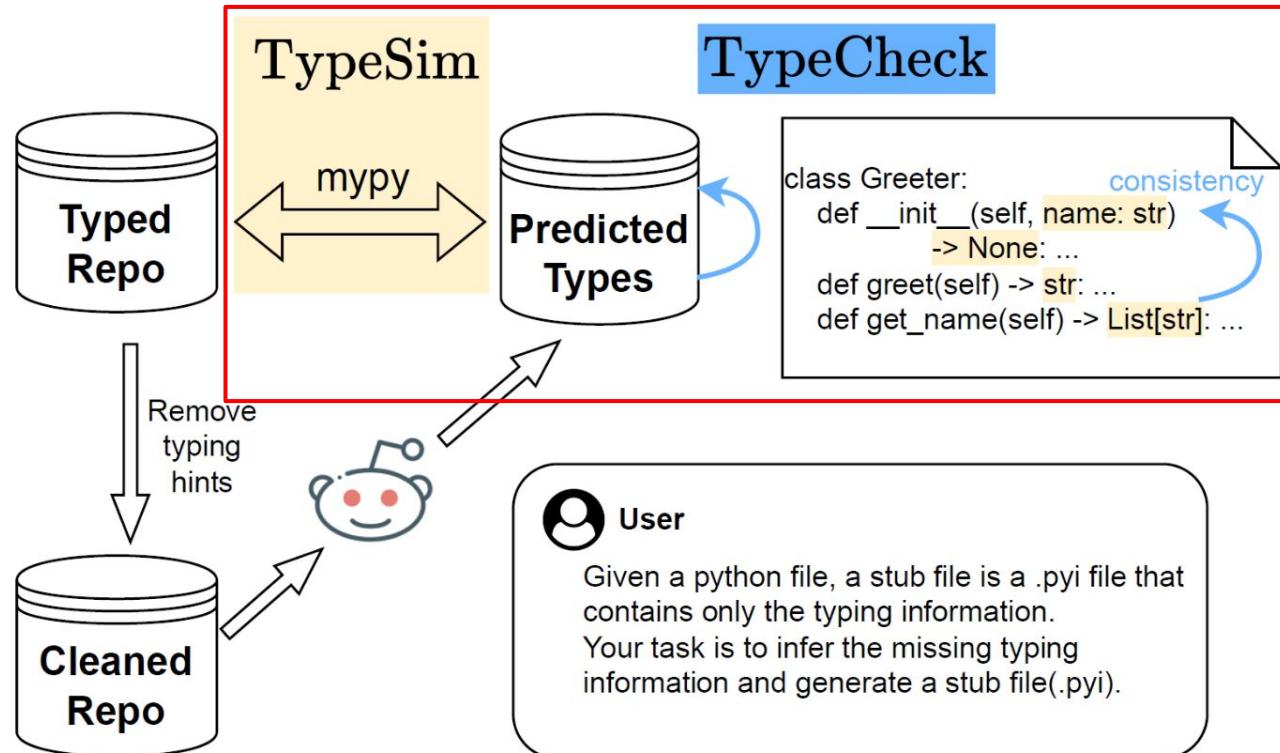
TypyBench

```
greeting.py
```

```
class Greeter:  
    def __init__(self, name:str)  
        -> None:  
        self.name = name  
    def greet(self) -> str:  
        return f"Hello, {self.name}!"  
    def get_name(self) -> str:  
        return self.name
```

```
greeting_without_types.py
```

```
class Greeter:  
    def __init__(self, name):  
        self.name = name  
    def greet(self):  
        return f"Hello, {self.name}!"  
    def get_name(self):  
        return self.name
```



Two Novel Metrics – TypeSim

functional similarity + structural relationships

Two Novel Metrics – TypeSim

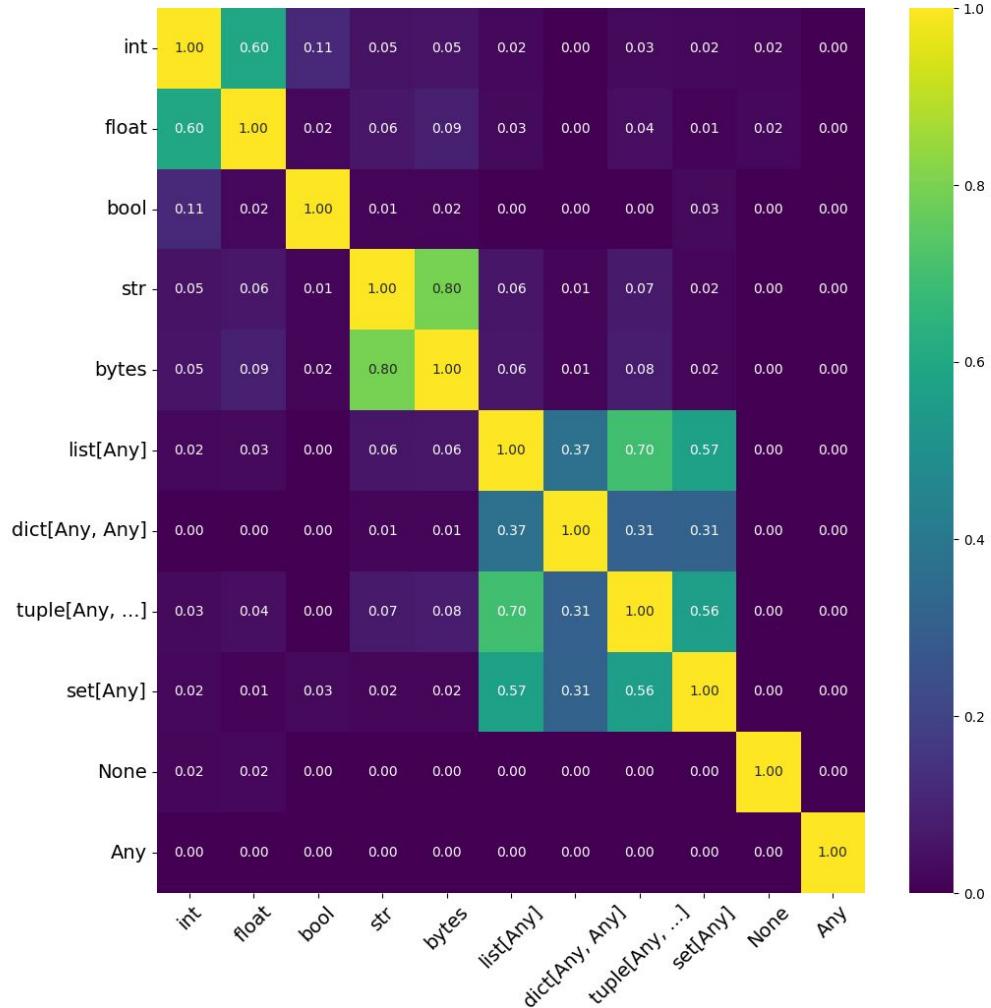
functional similarity + structural relationships

$$S(T1, T2) = \frac{|\text{attrs}(T1) \cap \text{attrs}(T2)|}{|\text{attrs}(T1) \cup \text{attrs}(T2)|}$$

Iterable vs Sequence: 0.92 > **List vs Sequence:** 0.7

int vs **float** : 0.6 > **int** vs **str** : 0.06

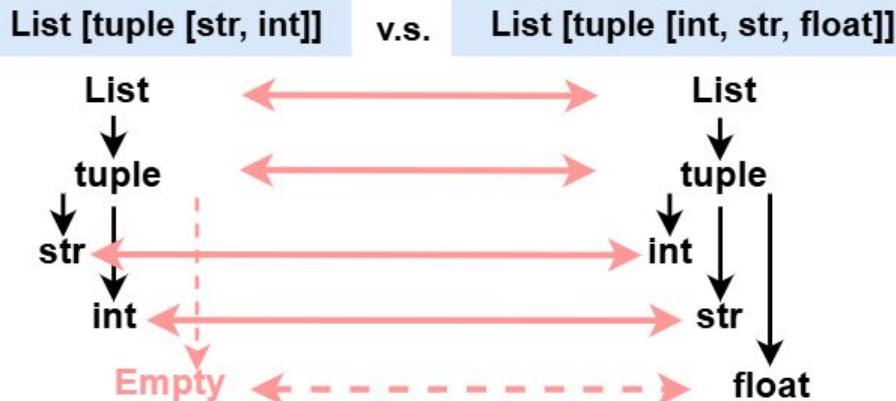
Similarity Score for Build-in Python Types



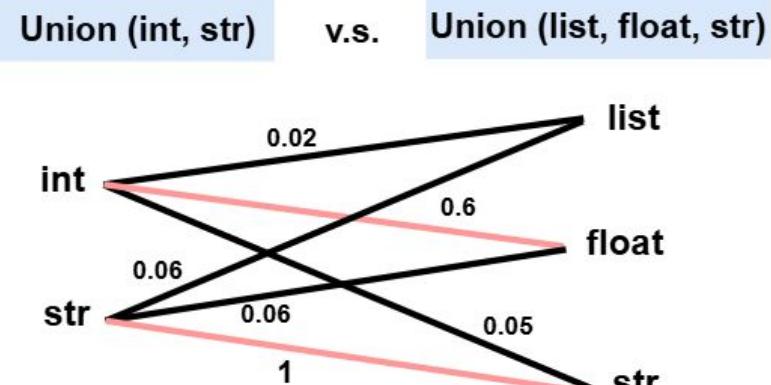
Two Novel Metrics – TypeSim

functional similarity + structural relationships

Ordered Arguments



Permutable Arguments



Two Novel Metrics – TypeCheck

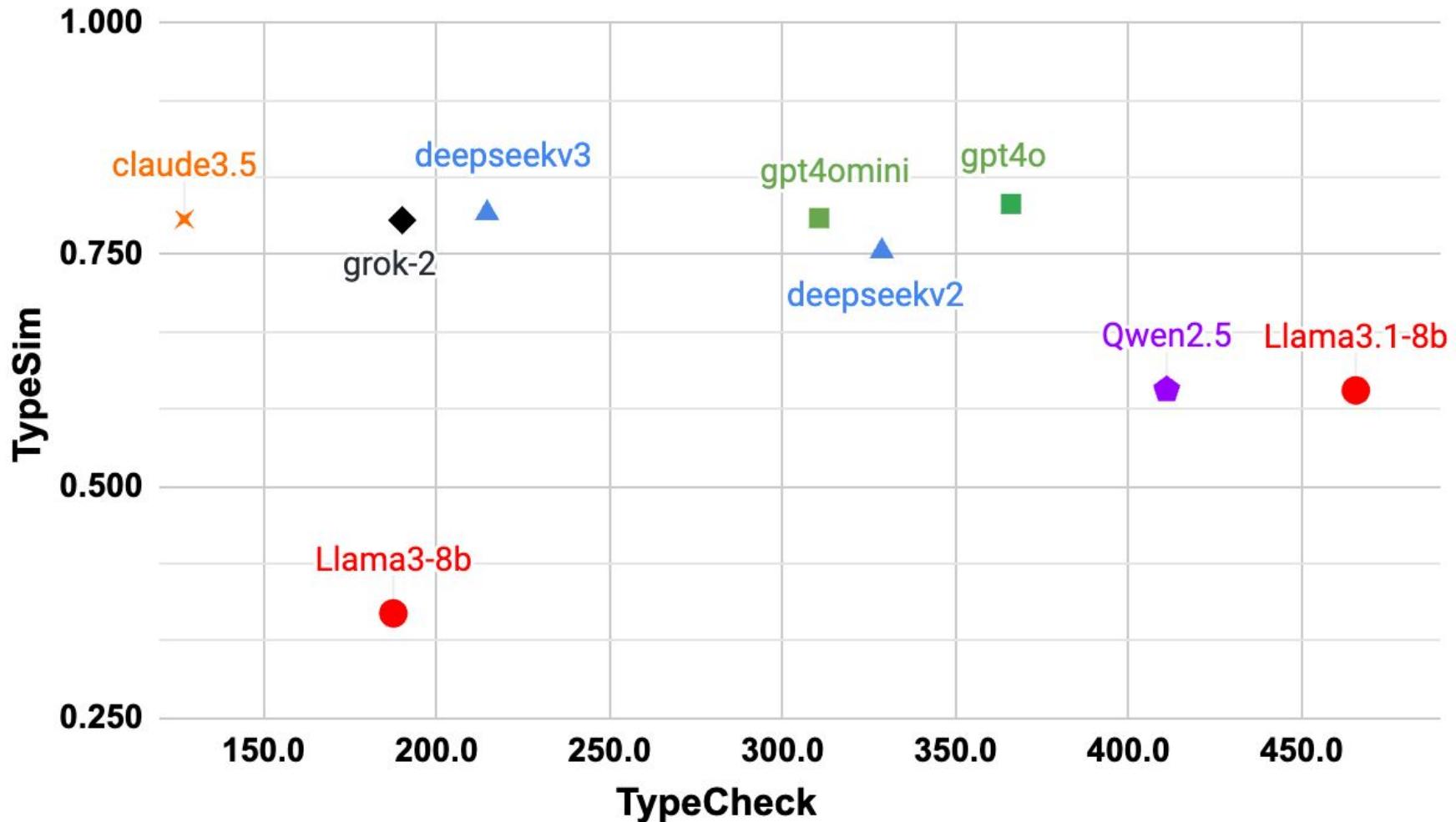
- Predicted types form a **coherent** system
- Errors generated by static type checker:
 - Incompatible return types
 - Invalid argument types

TypyBench Curation

We curate a benchmark dataset **TypyBench**, containing **50 well-typed Python repositories** from **GitHub** and **PyPI**.

Selection metrics: **Type Annotation Rate + Popularity + Complexity**

	# Repos	# Tokens	# Functions	# Cases
Train	20	8403760	31161	59966
Validation	10	4037666	13988	20177
Test	20	4983025	25101	46166



Takeaways

- **Complementary metrics**
 - **TYPESIM** captures **functionality + structural similarity**
 - **TYPECHECK** highlights **consistency** issues.

Takeaways

- **Complementary metrics**
- **Future directions:**
 - Repo-level focus to improve **cross-file type consistency**
 - **Context length trade-offs:** challenges in handling large inputs/outputs