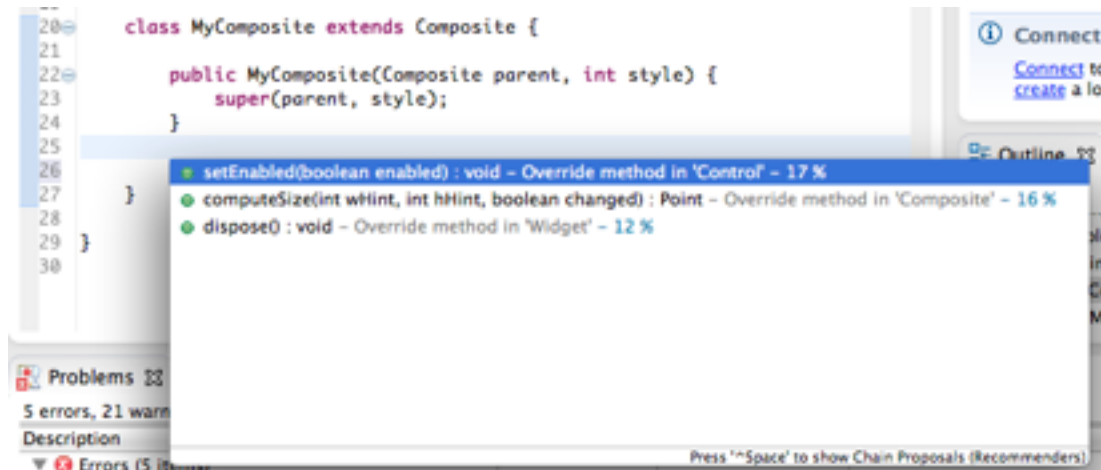# Statistical Models of Typed Syntax Trees

**Cyrus Omar**

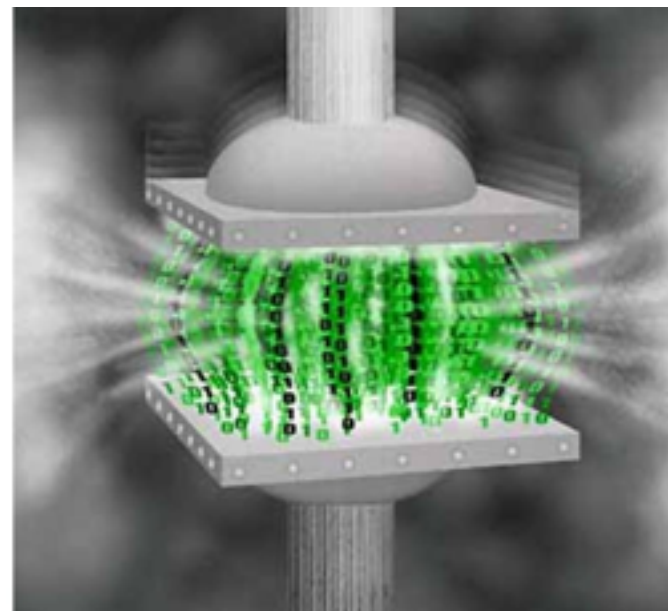**Carnegie Mellon University**

# Many tools could benefit from an understanding of the statistics of *natural programs*.
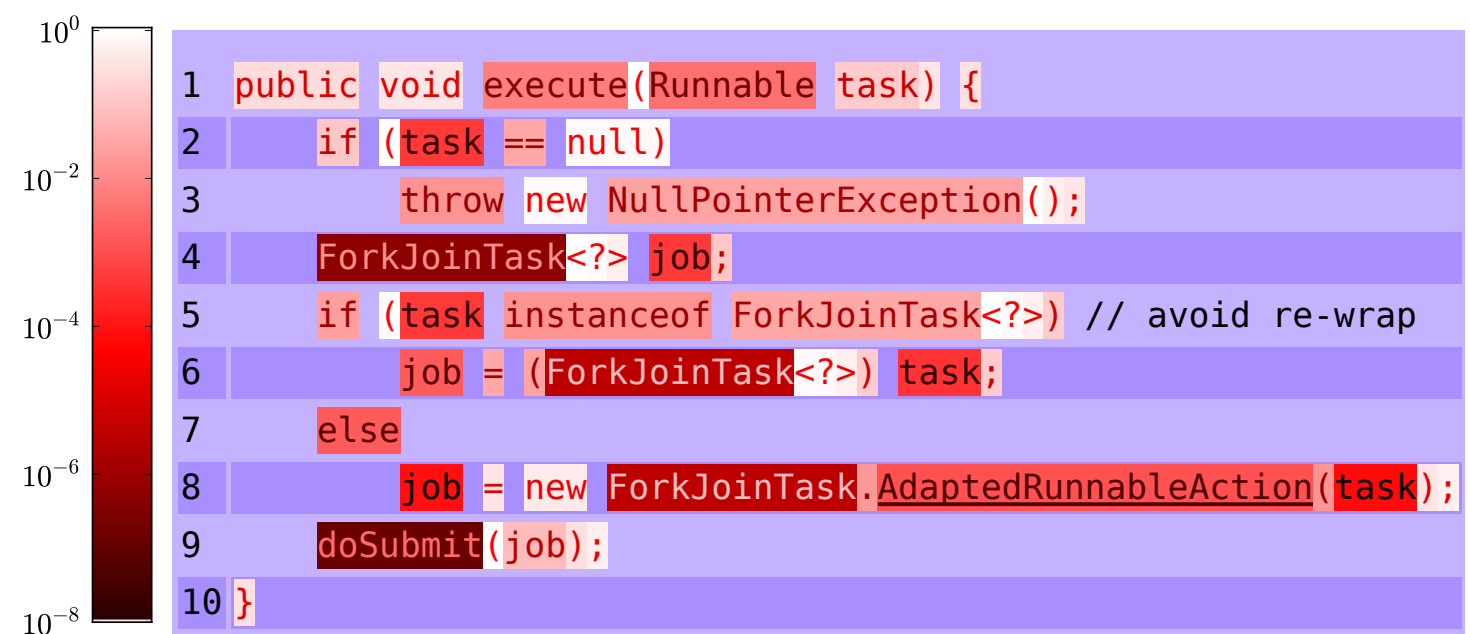


**code completion
engines**



**predictive programming
interfaces**



**code compression
engines**

Many tools could benefit from an understanding of the statistics of *natural programs*.



```
1  public void execute(Runnable task) {
2      if (task == null)
3          throw new NullPointerException();
4      ForkJoinTask<?> job;
5      if (task instanceof ForkJoinTask<?>) // avoid re-wrap
6          job = (ForkJoinTask<?>) task;
7      else
8          job = new ForkJoinTask.AdaptedRunnableAction(task);
9      doSubmit(job);
10 }
```

**code smell detectors**

```
enum Planet {
  MERCURY,
  VENUS,
  …,
  NEPTUNE
}

class Astronomer {
  void observe(Planet p) {…}
}
```

```
class Sagan extends Astronomer {
  void beforeBed() {
    |observe(favorite)|
  }

  Planet favorite = Planet.NEPTUNE;
}
```

To develop a probability distribution for code, we need to first choose a **representation of code**

```
enum Planet {
  MERCURY,
  VENUS,
  …,
  NEPTUNE
}

class Astronomer {
  void observe(Planet p) {…}
}
```

```
class Sagan extends Astronomer {
  void beforeBed() {
    observe(favorite)
  }

  Planet favorite = Planet.NEPTUNE;
}
```

$P($ observe $|$ ( favorite ) $|$ ( ) { $)$

- *n*-grams **(Hindle et al., ICSE 2012; Allamanis & Sutton, MSR 2013)**
- **topic modeling + part of speech analysis (Nguyen et al., FSE 2013)**

**Previous Work: Programs are Token Sequences**

```
enum Planet {
  MERCURY,
  VENUS,
  …,
  NEPTUNE
}

class Astronomer {
  void observe(Planet p) {…}
}
```

```
class Sagan extends Astronomer {
  void beforeBed() {
    observe(favorite)
  }

  Planet favorite = Planet.NEPTUNE;
}
```

P( call var this .observe( field var this .favorite ) )

$\mid \rho$ = **stmt,**

role

Our Approach: **Programs are Syntax Trees**

```
enum Planet {
  MERCURY,
  VENUS,
  …,
  NEPTUNE
}

class Astronomer {
  void observe(Planet p) {…}
}
```

```
class Sagan extends Astronomer {
  void beforeBed() {
    observe(favorite)
  }

  Planet favorite = Planet.NEPTUNE;
}
```

$\mathbf{P(}$ call var this .observe( field var this .favorite )

$| \rho = \mathbf{stmt}, \mathcal{T} = \mathbf{void}, \Gamma)$

role          type          typing
                            context

**Our Approach: Programs are *Typed* Syntax Trees**

```
enum Planet {
  MERCURY,
  VENUS,
  …,
  NEPTUNE
}

class Astronomer {
  void observe(Planet p) {…}
}
```

```
class Sagan extends Astronomer {
  void beforeBed() {
    observe(favorite)
  }

  Planet favorite = Planet.NEPTUNE;
}
```

$$= \mathbf{P}(\boxed{\text{call}\;\boxed{\text{var}\;\texttt{this}}\;.\texttt{observe(}\boxed{\text{field}\;\boxed{\text{var}\;\texttt{this}}\;.\texttt{favorite}})}$$

$$| \; \phi = \texttt{call}, \; \rho = \texttt{stmt}, \; \mathcal{T} = \texttt{void}, \; \Gamma) \, \mathbf{P}(\phi | \rho, \tau)$$

syntactic        role        type        typing
form                          context

**Bayes' rule!**

```
enum Planet {
  MERCURY,
  VENUS,
  …,
  NEPTUNE
}

class Astronomer {
  void observe(Planet p) {…}
}
```

```
class Sagan extends Astronomer {
  void beforeBed() {
    observe(favorite)
  }

  Planet favorite = Planet.NEPTUNE;
}
```
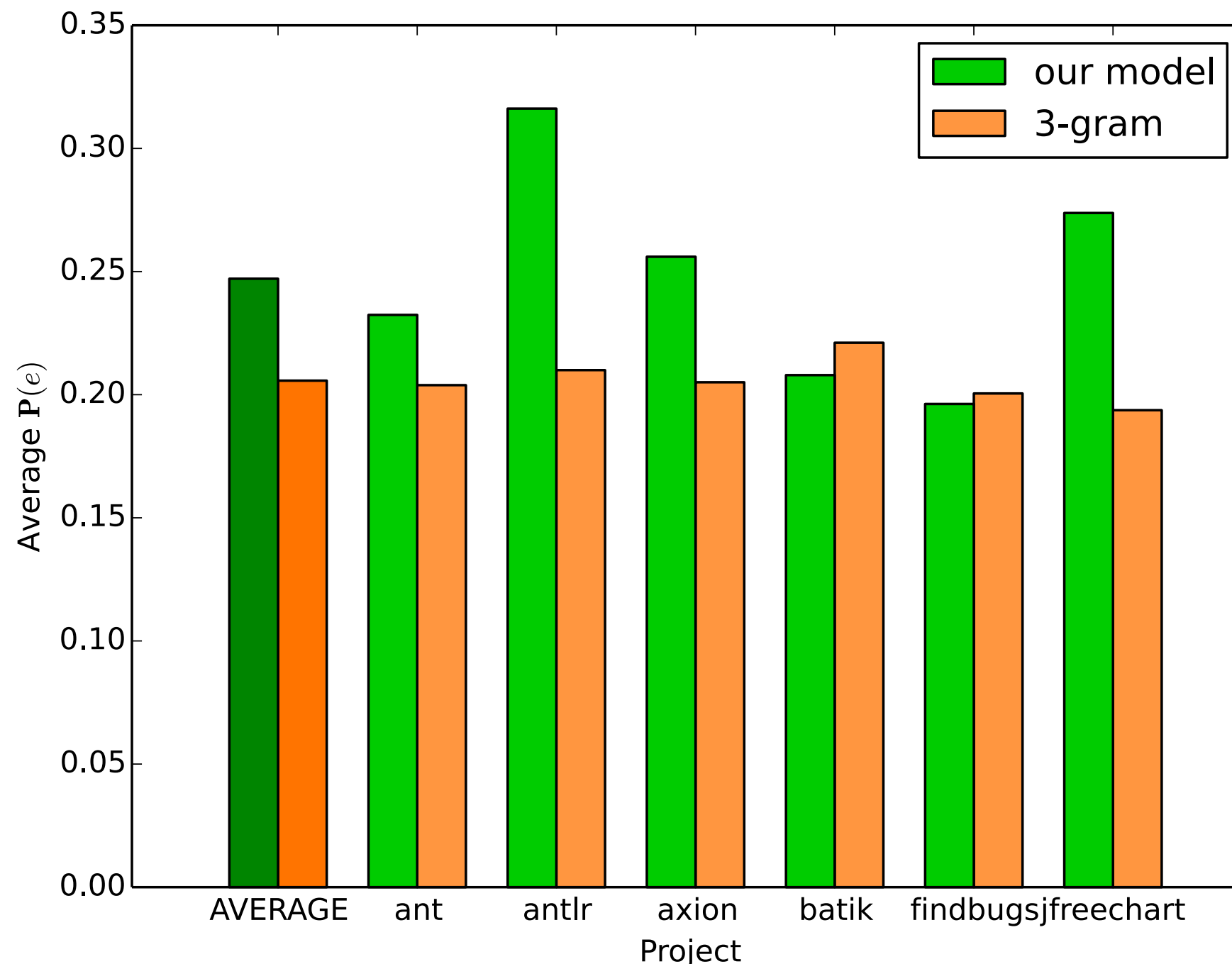
$$= \mathbf{P}(\boxed{\text{call } \boxed{\text{var this}} \text{.observe(} \boxed{\text{field } \boxed{\text{var this}} \text{.favorite}} )}$$

$$| \phi = \mathbf{call}, \rho = \mathbf{stmt}, \mathcal{T} = \mathbf{void}, \Gamma) \mathbf{P}(\phi | \rho, \tau)$$

$$= \mathbf{P}(\texttt{Sagan.observe} \mid \rho = \mathbf{stmt}, \mathcal{T} = \mathbf{void}, \Gamma)$$

$$\mathbf{P}(\boxed{\text{var this}} \mid \rho = \mathbf{targ}, \mathcal{T} = \mathbf{Sagan}, \Gamma)$$

$$\mathbf{P}(\boxed{\text{field } \boxed{\text{var this}} \text{.favorite}} \mid \rho = \mathbf{arg}, \mathcal{T} = \mathbf{Planet}, \Gamma)$$

$$\mathbf{P}(\phi | \rho, \tau)$$

We are starting to **implement** this model for Java using the **Eclipse JDT** for parsing and keeping track of $\Gamma$

http://www.github.com/cyrus-/syzygy

To test our implementation, we perform **10-fold cross-validation** of our model on a corpus of several large open source projects and **compare it to the 3-gram model** used in Hindle et. al, 2012.

We are taking a first-principles approach to
**source code prediction**
that combines the foundational techniques of both
**statistics** and **semantics**.

**P**( call var this .observe( field var this .favorite )

$| \rho =$ **stmt**, $\mathcal{T} =$ **void**, $\Gamma$)

role        type      typing
context