

Bayou

# Задача

- Генерация кода
  - Неполная спецификация
  - Вычислительная сложность

# Концепция

- Байесовский подход
  - Апостериорное распределение программ по данным evidence
- Скетчи
  - Language agnostic
  - Близки к intent программы
  - Простые
  - Могут быть конкретизированы в программу
- Скрытая переменная intent
  - Неявно моделируем глубинный смысл программы

# Workflow

- $X$  - evidence,  $Y$  - sketch,  $Z$  - intent
  - $P(Z|X)$  -- по evidence получаем распределение intent
  - Сэмплируем  $Z$
  - $P(Y|Z)$  -- по intent получаем распределение sketch
  - Сэмплируем  $Y$
  - Конкретизируем sketch в программу используя  $P(\text{Prog}|Y)$

\* На каждом этапе может применяться некоторое количество магии (и тервера)

# Пример

```
public class Program {  
    void foo(String file) {  
        Evidence.apical("readLine");  
    }  
}
```

```
import java.io.FileReader;  
import java.io.IOException;  
import java.io.FileNotFoundException;  
import java.io.BufferedReader;
```

```
public class Program {  
    void foo(String file) {  
        String s;  
        BufferedReader br;  
        FileReader fr;  
        try {  
            fr = new FileReader(file);  
            br = new BufferedReader(fr);  
            while ((s = br.readLine()) != null) {}  
            br.close();  
        } catch (FileNotFoundException _e) {}  
        } catch (IOException _e) {}  
    }  
}
```

# Постановка задачи

- $X$  -- evidences,  $Prog$  -- program,  $Env$  -- type environment
- Задача обучения: Имея корпус  $\{(X_i, Prog_i, Env_i)\}$  и модель  $P(X, Prog | M, Env)$  найти

$$M^* = \arg \max_M \sum_i \log P(X_i, Prog_i | M, Env_i)$$

- Задача синтеза: Имея evidence  $X$ , type env.  $Env$  и распределение  $P(X, Prog | Env)$  получить апостериорное распределение  $P(Prog | X, Env, safe(Prog, Env))$

# Скетчи

- $Y$  -- sketch
  - Скетчи определены в некоторой грамматике скетчей
  - Она не содержит переменных
  - Она содержит конструкции управления потоком выполнения
- Функция абстракции  $a: Prog \rightarrow Y$ 
  - Удаляет из AML\* всё, что не присуще грамматике скетчей
- Распределение конкретизации  $P(Prog | Y, Env)$ 
  - Если  $Y$  не конкретизируется под  $Env$ , то вернётся  $\perp$

\* AML -- процедурное подмножество Java

# Реформулировка задачи

- $X$  -- evidences,  $Prog$  -- program,  $Env$  -- type environment
- Задача обучения: Имея корпус  $\{(X_i, Prog_i, Env_i)\}$ ,  $a$  и  $P(X, Y | M)$  и  $Y_i = a(Prog_i)$  найти

$$M^* = \operatorname{argmax}_M \sum_Y \log P(X_i, Y_i | M)$$

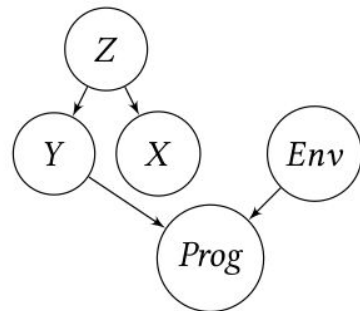
- Задача синтеза: Имея evidence  $X$ , type env.  $Env$  и распределение  $P(X, Y)$  получить апостериорное распределение

$$P(Prog | X, Env, safe(Prog, Env)) = \sum_Y P(Prog | Y, Env, Prog \neq \perp) P(Y | X).$$



# Intent переменная

- $Z$  -- скрытая переменная выражающая intent программы
  - $P(X_i, Y_i | M) = \int_{\Omega_z} P(X_i, Y_i, Z | M) dz$
- $P(X, Y, Z)$  тренируется на корпусе
- В работе используем  $P(Z|X)$ ,  $P(Y|Z)$



# BED

- $P(Z)$  -- нормальное многомерное распределение ( $\text{Normal}(O, I)$ )
- $P(X|Z)$  -- encoder
  - Кодлируем  $X_{\text{cxt}}, X_{\text{typ}}, X_{\text{api}}$  (используем  $f_{\text{cxt}}, f_{\text{typ}}, f_{\text{api}}$ )
  - $P(X|Z)$  -- сумма произведений нормальных распределений ( $\sum_{\text{ev}} \prod_i \text{Normal}_{\text{ev}}(Z, \sigma_{\text{ev}}^2 I)$ )
  - Отсюда выражается и  $P(Z|X)$
- $P(Y|Z)$  -- decoder
  - Используя RNN генерируем скетч
- Тренируется используя variational bound (max L, см. VAE)

$$\log P(X, Y | \mathbf{M}) = D_{KL}[P(Z | \mathbf{M}, X) || P(Z | \mathbf{M}, X, Y)] + \mathcal{L}$$

- Параметры --  $\sigma_{\text{ev}}, P(Y|Z)$  params, параметры функций кодирования

# Combinatorial concretization

- Сэмплируем программу из  $P(\text{Prog}|\Upsilon, \text{Env})$
- Но для простых программ можно обойтись random walk
  - Конкретизируем  $\Upsilon$  до некоторого промежуточного  $\Upsilon'$ 
    - Добавляем пару переменных
  - Находим возможные продолжения конкретизации
  - Случайно выбираем следующее

# Общая структура

- Evidence embedding layer
  - LDA преобразующий evidence в распределение над топиками
- Evidence encoder
  - NN преобразующая распределение в d-dim вектор
- Пространство intent переменной
- Intent decoder
  - Две RNN генерирующие скетчи по intent
- Combinatorial concretization
  - Случайный программы конкретизирующей скетч в текущем Env

# Эксперименты

- 98619 программ использующих Android API
  - Извлечены evidence
  - Построены скетчи
  - 80% обучение, 20% контроль
- 256, 64 and 32 топики для API calls, types, и contexts соответственно
- 128 узлов в encoder для API calls, 64 для types и contexts
- 256 узлов в decoder
- 32 -- размерность пространства intent
- 100 эпох

|             | Min | Max | Median | Mean |
|-------------|-----|-----|--------|------|
| $X_{Calls}$ | 1   | 9   | 2      | 2.6  |
| $X_{Types}$ | 1   | 7   | 2      | 1.8  |
| $X_{Cxt}$   | 0   | 11  | 2      | 2.1  |
| $X$         | 1   | 23  | 6      | 6.5  |

# Топики

| <i>Calls</i>   | <i>Types</i>  | <i>Cxt</i>  |
|--|---|---|
| 1. setJavaScriptEnabled<br>setWebViewClient<br>setWebChromeClient      | 1. Builder<br>AlertDialog<br>Dialog                 | 1. Reader<br>InputStream<br>String                        |
| 2. getRunningAppProcesses<br>restartPackage<br>killBackgroundProcesses | 2. InputStreamReader<br>BufferedReader<br>Throwable | 2. Notification<br>Intent<br>Context                      |
| 3. setTitle<br>setMessage<br>show                                      | 3. LayoutParams<br>LinearLayout<br>FrameLayout      | 3. OnItemClickListener<br>OnScrollListener<br>ListAdapter |

# Quality scenario

- Запись в файл int
  - API call “write”, type “FileWriter”, context “int “
- InputStream из BluetoothSocket
  - context “BluetoothAdapter”, “String”