# Tutorial of Bee

## 1 Introduction

Bee is a program synthesis framework that builds on the **data-actions** model. That is, a PBE task can be viewed as processing some **data entities** with repetitive **user actions**. For example, in a task to download hundreds of attachments from an email client, the data are the emails and attachments, and the user actions are the downloading operations on some attachments.

If the target PBE tasks match this model, developers can provide specifications of the data and actions to customize a PBE synthesizer on Bee.

Specifically, developers need to represent their data as relational tables so that Bee can correctly process them. Also, developers need to specify the user actions as program functions/methods. In this way, Bee can record these methods' invocations as examples, and invoke these methods after the target actions are generated.

In the followings, we would use a PBE task as running example to introduce how to create a synthesizer appropriate to the given tasks with Bee step by step (Section 2 to 5). We then briefly introduce the internals of Bee (Section 6) and discuss how to design more effective table schemas (Section 7).

## 2 Developing PBE for Email Client

Now imagine you are a developer of an email client and you want to support an auto-download function. To start with, consider the following specific tasks from the end-users of your email client.

**Task A**: From emails named with "assignment", save attachments to folders named by students' email addresses. Figure 1 shows three example emails. Two emails are from students and the three attachments included should be saved to corresponding target paths. For other irrelevant emails, no actions should be taken.

| Email Title | Sender | Attachment | Target Path |
|:---:|:---:|:---:|:---:|
| assignment1 | david@xxx.edu | project.zip | ~/downloads/A1/david |
| | | report.pdf | ~/downloads/A1/david |
| assignment1 | lucy@xxx.edu | assign1.zip | ~/downloads/A1/lucy |
| Discount Now | adv@yyy.com | discount.jpg | N/A |

Figure 1: Examples for task A.

**Task B**: Download bills sent from a bank to folders named by year-months. Figure 2 shows four example emails. Three emails are from the bank and the attached bills are saved to corresponding target paths. One email is irrelevant and no action is taken.

| Email Title | Sender | Date | Attachment | Target Path |
|:---:|:---:|:---:|:---:|:---:|
| Billing info | bill@boc.com | 2021-01-31 | bill.pdf | ~/bills/2021/Jan |
| Billing info | bill@boc.com | 2021-02-28 | bill.pdf | ~/bills/2021/Feb |
| Billing info | bill@boc.com | 2021-03-31 | bill.pdf | ~/bills/2021/Mar |
| Discount Now | adv@yyy.com | 2021-03-31 | discount.jpg | N/A |

Figure 2: Examples for task B.

Task A and Task B are two typical scenarios that can be automated by PBE. As a developer of the email client, you would like to add such a PBE feature to your software. With this PBE feature, not only Task A and Task B, but also other similar user tasks can be benefited. You may consider extending the PBE features for more scenarios in practice, but in this tutorial let's just focus on these two tasks.

Compared with developing PBE from scratch, using program synthesis framework like BEE is more convenient. Specifically, we need to supply BEE with the table representation of data (Section 3) and methods corresponding to the user actions (Section 4).

## 3   Modeling Data

We now consider what are the data in email client and design table schemas to represent them. An email has fields such as title, content, sender, receivers, datetime, attachments, to which email it is replied, and so on. These

fields can all be modeled in relational tables, but in this example PBE we only model some of them. Readers can follow the approach and apply to other fields.

Since an email may include multiple attachments, we consider modeling attachments in a separate table. Figure 3 shows the table modeling of data in Figure 1. The first table holds basic email information (`string` type) plus an `Eid` to mark the identity of each email (`id` type). An `id` value can be understood as the address/pointer/reference of an object. The second table holds attachments' information. Each row (attachment) has a filename, a distinct `Fid` and the `Eid` of the email it belongs to.

| Eid | Title | Year | Month | Sender |
|-----|-------|------|-------|--------|
| e1 | assignment1 | 2021 | Mar | david@xxx.edu |
| e2 | assignment1 | 2021 | Mar | lucy@xxx.edu |
| e3 | Discount Now | 2021 | Mar | adv@yyy.com |

| Eid | Fid | Filename |
|-----|-----|----------|
| e1 | f1 | project.zip |
| e1 | f2 | report.pdf |
| e2 | f3 | assign1.zip |
| e3 | f4 | discount.jpg |

Figure 3: Relational-table modeling of emails and attachments in Figure 1.

## 3.1 APIs for modeling data

To implement the above design, BEE provides a set of annotation-based APIs. For each table schema, you need to define a class, then mark it with the `[Entity]` annotation (using C#). If you are using Java or Python, the annotation should be `@Entity`. Below shows the declarations of two table schemas in C#.

```
1 [Entity]
2 public class EmailEntity {
3     ...
4 }
5
6 [Entity]
7 public class AttachmentEntity {
8     ...
9 }
```

Then, for each column in a table schema, you need to define a non-static method with no parameters, and mark the method with `[Field]`. The type of column is the same as the return type of the method. For example, the codes below composites an `Email` object and define the four string columns in the email table.

```
1  [Entity]
2  public class EmailEntity {
3      private Email _email;
4      EmailEntity(Email email) {
5          _email = email;
6      }
7
8      [Field]
9      public string Title() {
10         return _email.Title;
11     }
12
13     [Field]
14     public string Sender() {
15         return _email.Sender;
16     }
17
18     [Field]
19     public string Year() {
20         return _email.SendTime.Year.ToString();
21     }
22
23     [Field]
24     public string Month() {
25         return _email.SendTime.Month.ToString();
26     }
27
28     ...
29 }
```

Regarding columns of `id` type, just define a method returning an object and mark the method with `[IdField]`. For example, the `Eid` method below returns an `Email` object, then the address/reference/pointer of this object would appear in the `Eid` column of this table.

```
1  [Entity]
2  public class EmailEntity {
3      ...
4
5      [IdField]
```

```
6      public Email Eid() {
7          return _email;
8      }
9  }
```

# 4  Actions

In addition to the data, we also need to specify the user actions as program functions/methods, including the method names, parameters, and implementations. For example, to specify a download operation, we may just name the method as `Download`, and have two parameters (the attachment to download, and where it is downloaded to) like `Download(Attachment attachment, string targetPath)`. This implementation, by definition, should save the `attachment` to `targetPath`.

The action instances for the data in Figure 3 are shown in Figure 4. Note that there is a one-to-one correspondence between a set of actions, and their table representation.

| action | attachment | target path |
|----------|------------|---------------------|
| download | f1 | ~/downloads/A1/david |
| download | f2 | ~/downloads/A1/david |
| download | f3 | ~/downloads/A1/lucy |

Figure 4: Example actions

## 4.1  APIs for action signature

To specify an action in C#, simply define a corresponding static method and mark it with `[Action]`. The codes below show the definition of download.

```
1  [Entity]
2  public class AttachmentEntity {
3      ...
4
5      [Action]
6      public static void Download(Attachment attachment,
7                                  string targetPath) {
8          attachment.SaveTo(targetPath);
9      }
10 }
```

5

# 5   Connect Everything

After defining the table schemas and action signatures, let us see how to connect them to the framework and invoke the synthesis procedure. The codes below demonstrate how to supply example data/actions and do synthesis.

```
 1 // prepared Email objects e1, e2, e3
 2 // and Attachment objects f1, f2, f3, f4
 3 ...
 4
 5 var emailsTable = new DataParser<EmailEntity>()
 6     .ParseAsTable(new List<EmailEntity> {
 7         new EmailEntity(e1),
 8         new EmailEntity(e2),
 9         new EmailEntity(e3)
10 });
11 var attachmentsTable = new DataParser<AttachmentEntity>()
12     .ParseAsTable(new List<AttachmentEntity> {
13         new AttachmentEntity(e1, f1),
14         new AttachmentEntity(e1, f2),
15         new AttachmentEntity(e2, f3),
16         new AttachmentEntity(e3, f4),
17 });
18
19 var recorder = ActionRecorder.GetInstance();
20 recorder.Start();
21 AttachmentEntity.Download(f1, "~/downloads/A1/david");
22 AttachmentEntity.Download(f2, "~/downloads/A1/david");
23 AttachmentEntity.Download(f3, "~/downloads/A1/lucy");
24 recorder.End();
25 var exampleActions = recorder.RecordedActions;
26
27 var synthesizer = new Synthesizer(new List<Table> {
28                                     emailsTable,
29                                     attachmentsTable },
30                                 exampleActions);
31 var program = synthesizer.Synthesize(10000).First();
```

# 6 Internals of Bee

You may wonder how BEE works. With the data represented as relational tables and actions' signatures provided, BEE ships a language extended from SQL to synthesize programs that generate actions from data tables. A BEE program first takes as input the tables converted from the data, then performs transformations such as filtering rows, aggregating information (max, count, etc.) along columns and so on. In principle, developers do not have to bother handling these transformations, but should understand what transformations are available (Section 6.2) in order to adjust the table schema for representing data.

After transformation, the BEE program would map appropriate columns in the transformed tables to actions' arguments, and generate corresponding action instances. Finally, the generated actions can be used to replace manual repetitive work.

## 6.1 Synthesized Program

In order to transform the tables in Figure 3 to the table in Figure 4, the program should be able to: (1) filter assignment emails by title; (2) join the two tables on the two `Eid` columns; (3) select columns to map to the parameters of the download signature, where the destination is by performing string transformation on the sender column.

Listing 1 shows the synthesized program for the above PBE task. The three statements correspond to the above steps.

```
1 function(t_emails, t_attachs)
2 {
3     t1 = Filter(t_emails, strContains(Title, "assignment"));
4     t2 = Join(t_attachs, t1, t_attachs.Eid, t1.Eid);
5     Yield(Download, t2, Fid, Mutate(concat(...), Sender));
6 }
```

Listing 1: Synthesized Program.

## 6.2 Grammar

Figure 5 shows the simplified grammar of the DSL in BEE. While developers do not need to modify this DSL, they can check the DSL's mechanism and adjust the data modeling and action signatures accordingly.

Overall, a program $p$ contains two consecutive subroutines, a transform program $p_t$ for transforming intermediate tables, and a mapping program

7

| | | | |
|---:|:---:|:---:|:---|
| Program | $p$ | $\coloneqq$ | $p_t\ p_m$ |

| | | | |
|---:|:---:|:---:|:---|
| Transform Program | $p_t$ | $\coloneqq$ | $s_t; p_t \mid \epsilon$ |
| Transform Statement | $s_t$ | $\coloneqq$ | $t = \text{Filter}\ (t,\ \phi)$ |
| | | $\mid$ | $t = \text{Join}\ (t_1,\ t_2,\ col_1^{\texttt{Id}},\ col_2^{\texttt{Id}})$ |
| | | $\mid$ | $t = \text{GroupJoin}\ (t,\ col_{index},\ (\alpha,\ col)\ldots)$ |
| | | $\mid$ | $t = \text{Order}\ (t,\ col,\ c_{start},\ c_{inv})$ |
| | | $\mid$ | $t = \text{Order}\ (t,\ col,\ c_{start},\ c_{inv},\ col_{index})$ |
| Aggregation | $\alpha$ | $\coloneqq$ | $\max \mid \min \mid \text{sum} \mid \text{avg} \mid \text{cnt}$ |
| Predicate | $\phi$ | $\coloneqq$ | $(\phi_1 \wedge \phi_2) \mid (\phi_1 \vee \phi_2) \mid \neg\phi$ |
| | | $\mid$ | $ps\ (col_1, col_2) \mid ps\ (col, c)$ |

| | | | |
|---:|:---:|:---:|:---|
| Mapping Program | $p_m$ | $\coloneqq$ | $s_m; p_m \mid \epsilon$ |
| Mapping Statement | $s_m$ | $\coloneqq$ | $\text{Yield}\ (c_{action}, t, \rho\ldots)$ |
| Projection | $\rho$ | $\coloneqq$ | $col \mid c \mid \text{Mutate}\ (f,\ col\ldots)$ |

Figure 5: Grammar of Bee. $\epsilon$ refers to the empty literal, $f$ refers to a feature, $t$ refers to a table variable, $col$ refers to a column name, $c$ refers to a constant value, $ps$ refers to a predicate symbol

$p_m$ for mapping table rows to the action signatures. The Filter, Join and GroupJoin transform statements correspond to the where, join and group clauses in SQL, respectively. In addition, the order statement can add columns indicating the orders of the rows.

The mapping statement (Yield) can map either a column or a constant value or a mutation of some columns to the parameter of an action signature. The Mutate operator can generate columns by calculating integral arithmetics and string concatenation on existing columns.

# 7 Design Guidelines

Considering that there could be multiple table schemas for representing the same data, and that the principles in SQL modeling (for persisting data) may not be appropriate in BEE (for program synthesis), here we provide a few guidelines on how to design effective table schemas for BEE.

*Redundancy is not necessarily harmful.* Consider the table for emails in Figure 3. Note that the sender column do not meet the principle to reduce redundancy in SQL. In SQL, considering that a same sender could appear multiple times, we should use a separate table to store the addresses and use keys to connect them to the email table. However, the table nesting the sender column is valid in BEE, since data redundancy usually does not harm to synthesizing programs in BEE. Instead, embedding them in one table would save the step of joining tables in some cases, and thus is more efficient for program synthesis.

*You can represent the same data in multiple schemas.* For example, we use two schemas in our spreadsheet PBE tool built atop BEE: (a) a spreadsheet is divided into cells, where each cell corresponds to one row in the relational table, with columns such as text content, row number and row heading; (b) the original tabular form, where rows in spreadsheet are mapped to rows in relational table and so for columns. Different table schemas could affect the performance of synthesis procedure. For example, if we want to project all the cell texts to one column, schema (b) requires extra transformation because the cell texts exist in different columns, while schema (a) can be used as is. In other cases, schema (b) may outperform schema (a). BEE allows the coexistence of both schemas and would automatically find the one with higher efficiency in the synthesis process, so developers can start with any valid schema, then provide some other schemas to optimize the performance for special scenarios.