# Clone Detection

John Businge

([john.businge@uantwerpen.be](mailto:john.businge@uantwerpen.be))
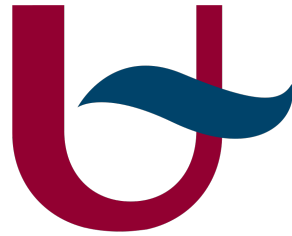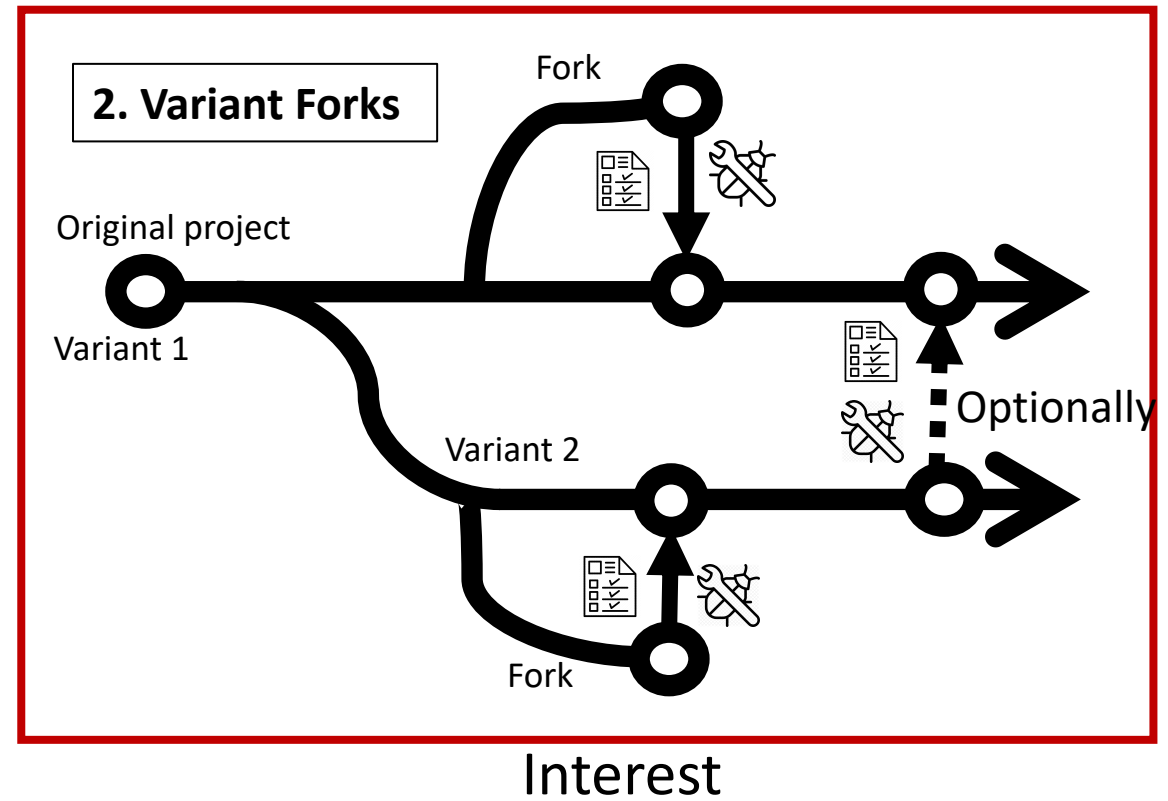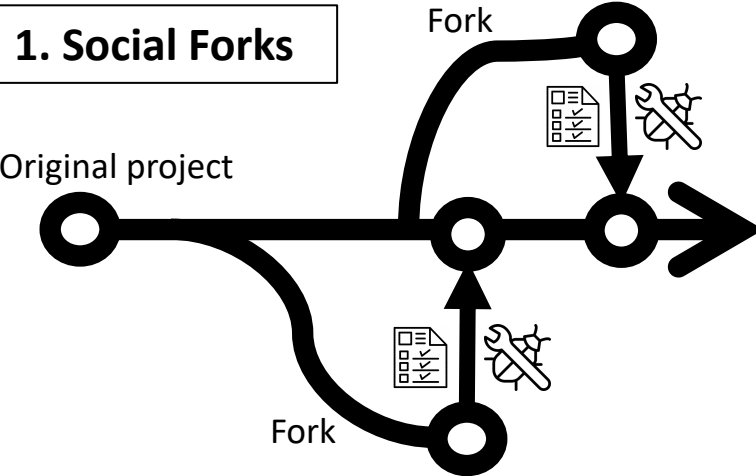
Universiteit
Antwerpen

# When is code duplication okay?

- If duplicating something hits a deadline and not duplicating doesn't, then I would rather deliver today and fix the tech debt tomorrow.

- There are two development paradigms
  - Clone-and-own (forking) - a new variant of a software system is created by copying and adapting an existing variant
  - Software product line - which consists of a set of similar software products with well-defined commonalities and variabilities.
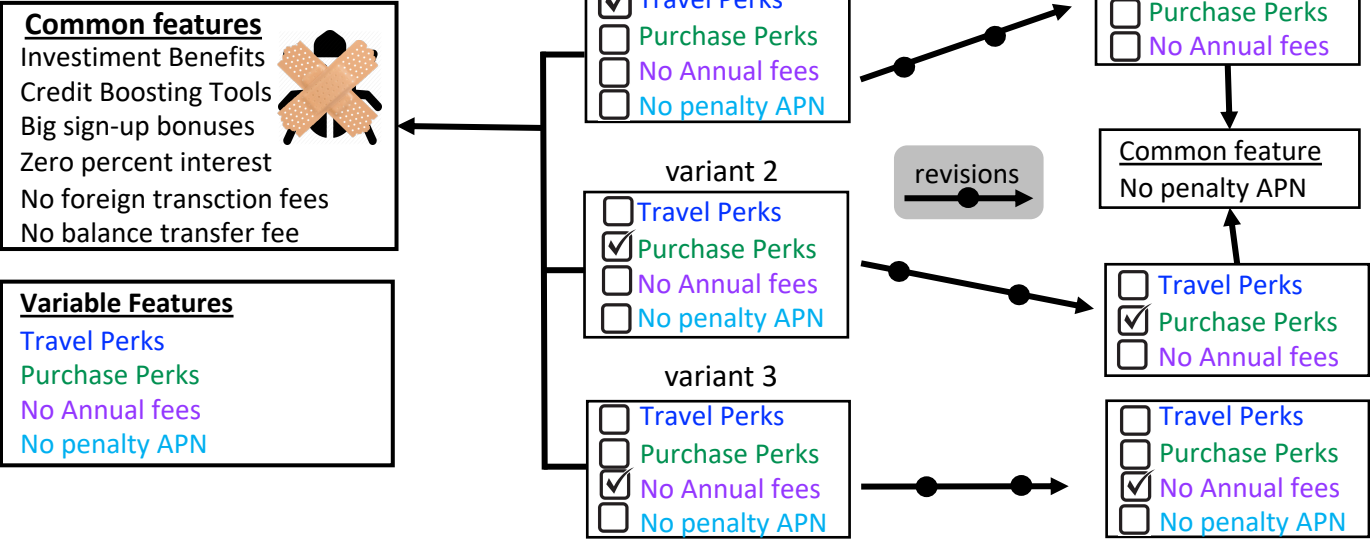
# Clone-and-own



1. Social Forks

Fork

Original project

Fork

2. Variant Forks

Fork

Original project

Variant 1

Variant 2

Optionally

Fork

Interest

# Engineering of multi-variant software systems (credit card variants)

## Software Product Line

**Common features**
Investiment Benefits
Credit Boosting Tools
Big sign-up bonuses
Zero percent interest
No foreign transction fees
No balance transfer fee

**Variable Features**
Travel Perks
Purchase Perks
No Annual fees
No penalty APN

variant 1
- ☑ Travel Perks
- ☐ Purchase Perks
- ☐ No Annual fees
- ☐ No penalty APN

variant 2
- ☐ Travel Perks
- ☑ Purchase Perks
- ☐ No Annual fees
- ☐ No penalty APN

variant 3
- ☐ Travel Perks
- ☐ Purchase Perks
- ☑ No Annual fees
- ☐ No penalty APN

revisions

- ☑ Travel Perks
- ☐ Purchase Perks
- ☐ No Annual fees

Common feature
No penalty APN

- ☐ Travel Perks
- ☑ Purchase Perks
- ☐ No Annual fees

- ☐ Travel Perks
- ☐ Purchase Perks
- ☑ No Annual fees
- ☐ No penalty APN

## Clone&own (forking/brancing)

v3
v2
v1
v4

- ☑ **Common**
- ☑ Travel Perks
- ☑ Purchase Perks
- ☑ No Annual fees

- ☑ **Common**
- ☐ Travel Perks
- ☑ Purchase Perks
- ☐ No Annual fees

- ☑ **Common**
- ☐ Travel Perks
- ☐ Purchase Perks
- ☐ No Annual fees

- ☑ **Common**
- ☑ Travel Perks
- ☐ Purchase Perks
- ☐ No Annual fees

## Benefites

- Enforces systemmatic reuse
- Easy to fix bugs
- Scales easily (variants)
- Common developers (easy to coodinate)
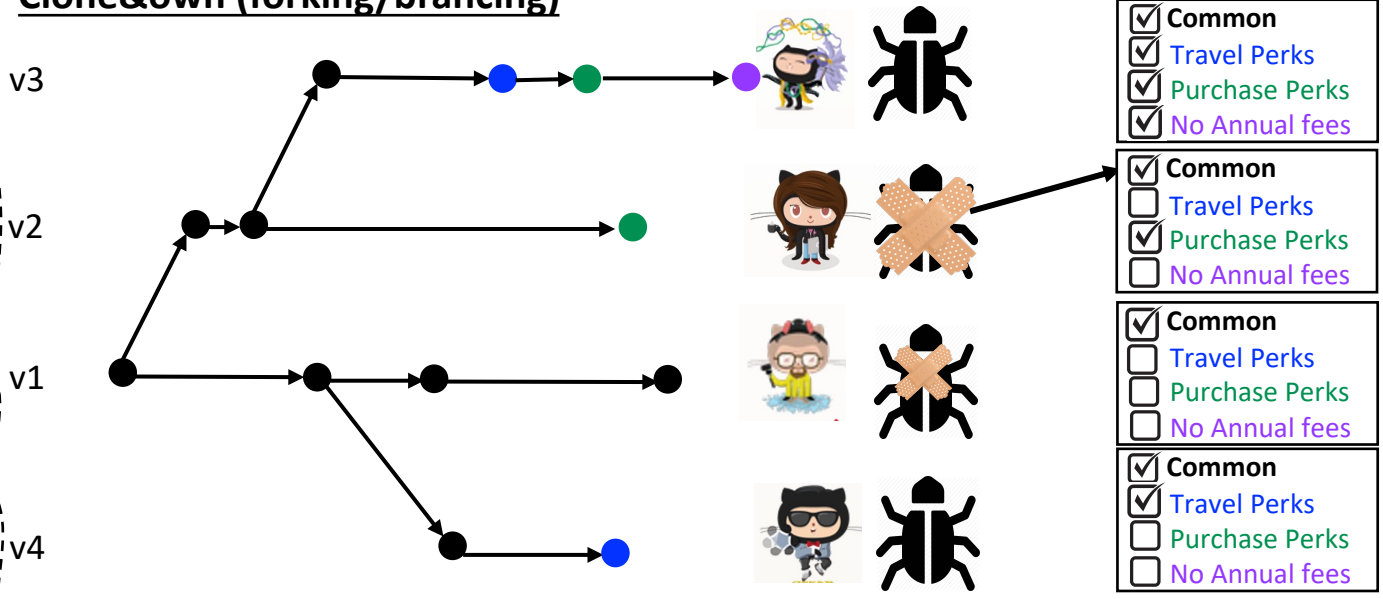- Saves long-term costs, time and effort (easy to manage)

## Limitations

- Upfront heavy investment
- Limited developer independency
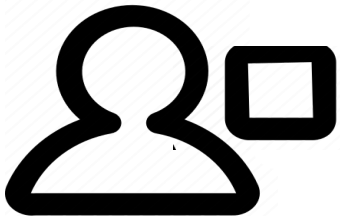- Not worth with few vatiants

## Benefits

- Flexibility in variant creation
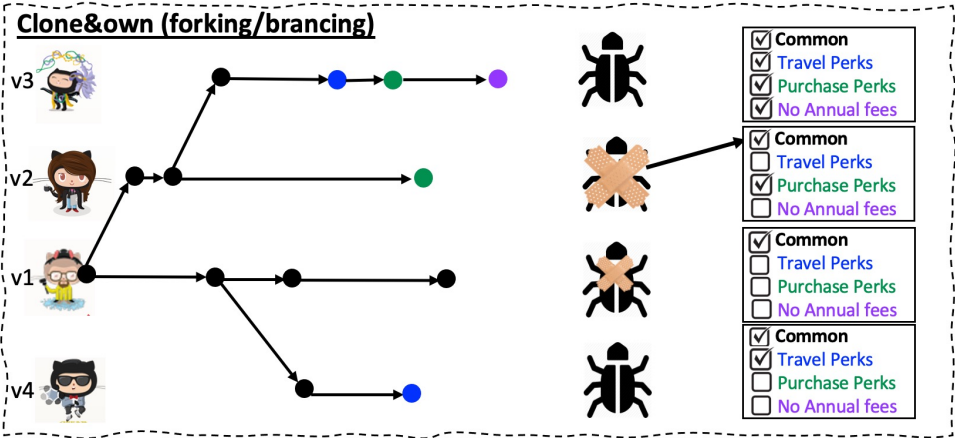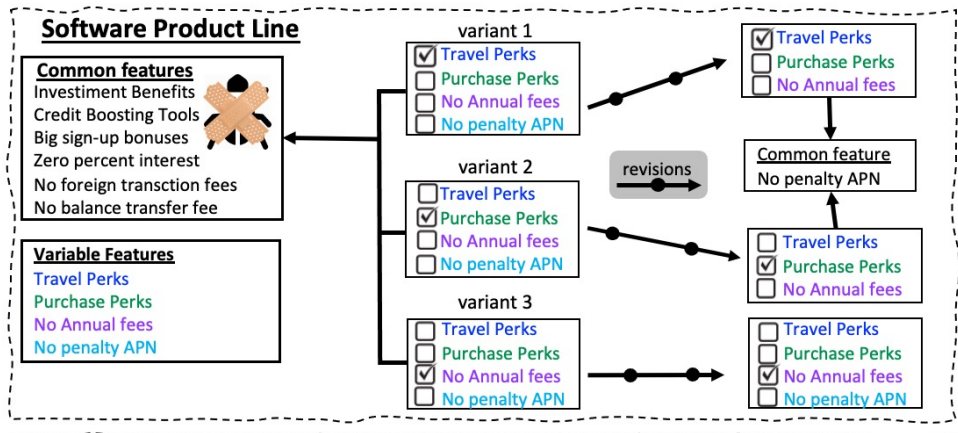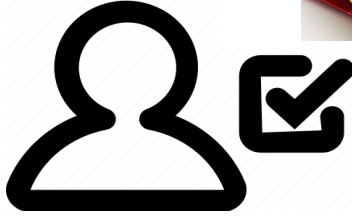- Limited initial costs
- Developer independency

## Limitations

- Redudancy/effort duplication
- Does not scale
- Misses opportunity (missed patches)
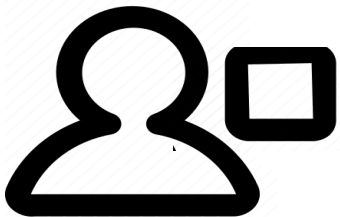- Diverse developers (difficult to coodinate)
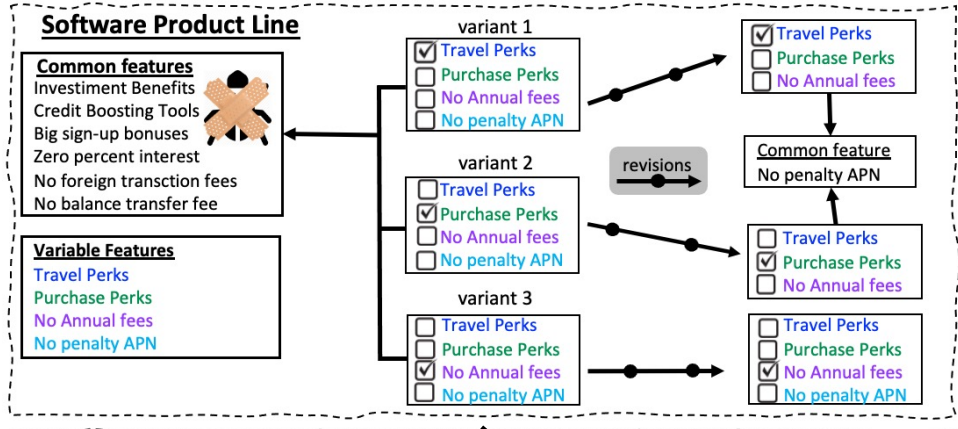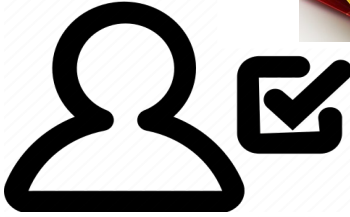- Expensive to engineer (SPL)

github SOCIAL CODING

Bitbucket

# State-of-practice



**Software Product Line**

**Common features**
Investiment Benefits
Credit Boosting Tools
Big sign-up bonuses
Zero percent interest
No foreign transction fees
No balance transfer fee

**Variable Features**
Travel Perks
Purchase Perks
No Annual fees
No penalty APN

variant 1
☑ Travel Perks
☐ Purchase Perks
☐ No Annual fees
☐ No penalty APN

variant 2
☐ Travel Perks
☑ Purchase Perks
☐ No Annual fees
☐ No penalty APN

variant 3
☐ Travel Perks
☐ Purchase Perks
☑ No Annual fees
☐ No penalty APN

revisions

☑ Travel Perks
☐ Purchase Perks
☐ No Annual fees

Common feature
No penalty APN

☐ Travel Perks
☑ Purchase Perks
☐ No Annual fees

☐ Travel Perks
☐ Purchase Perks
☑ No Annual fees
☐ No penalty APN

**Clone&own (forking/brancing)**

v3
v2
v1
v4

☑ Common
☑ Travel Perks
☑ Purchase Perks
☑ No Annual fees

☑ Common
☐ Travel Perks
☑ Purchase Perks
☐ No Annual fees

☑ Common
☐ Travel Perks
☐ Purchase Perks
☐ No Annual fees

☑ Common
☑ Travel Perks
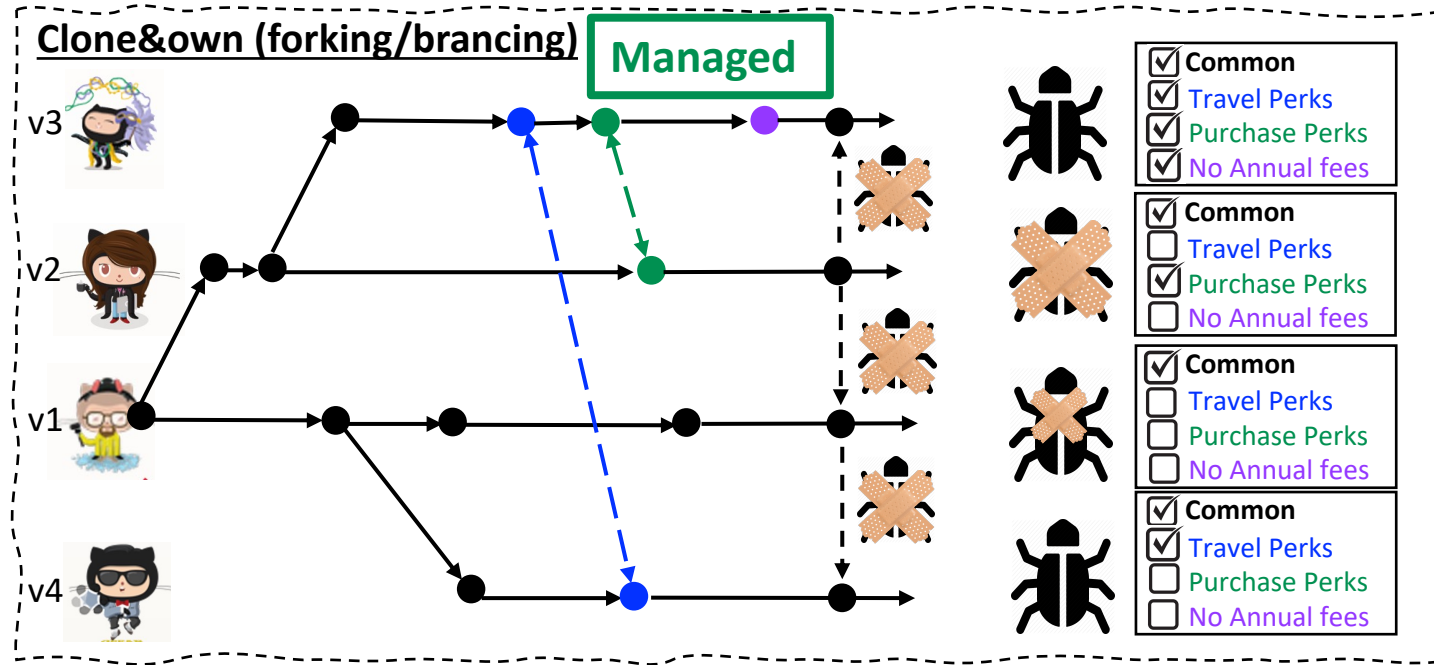☑ Purchase Perks
☐ No Annual fees

## Illustration of diverged variants
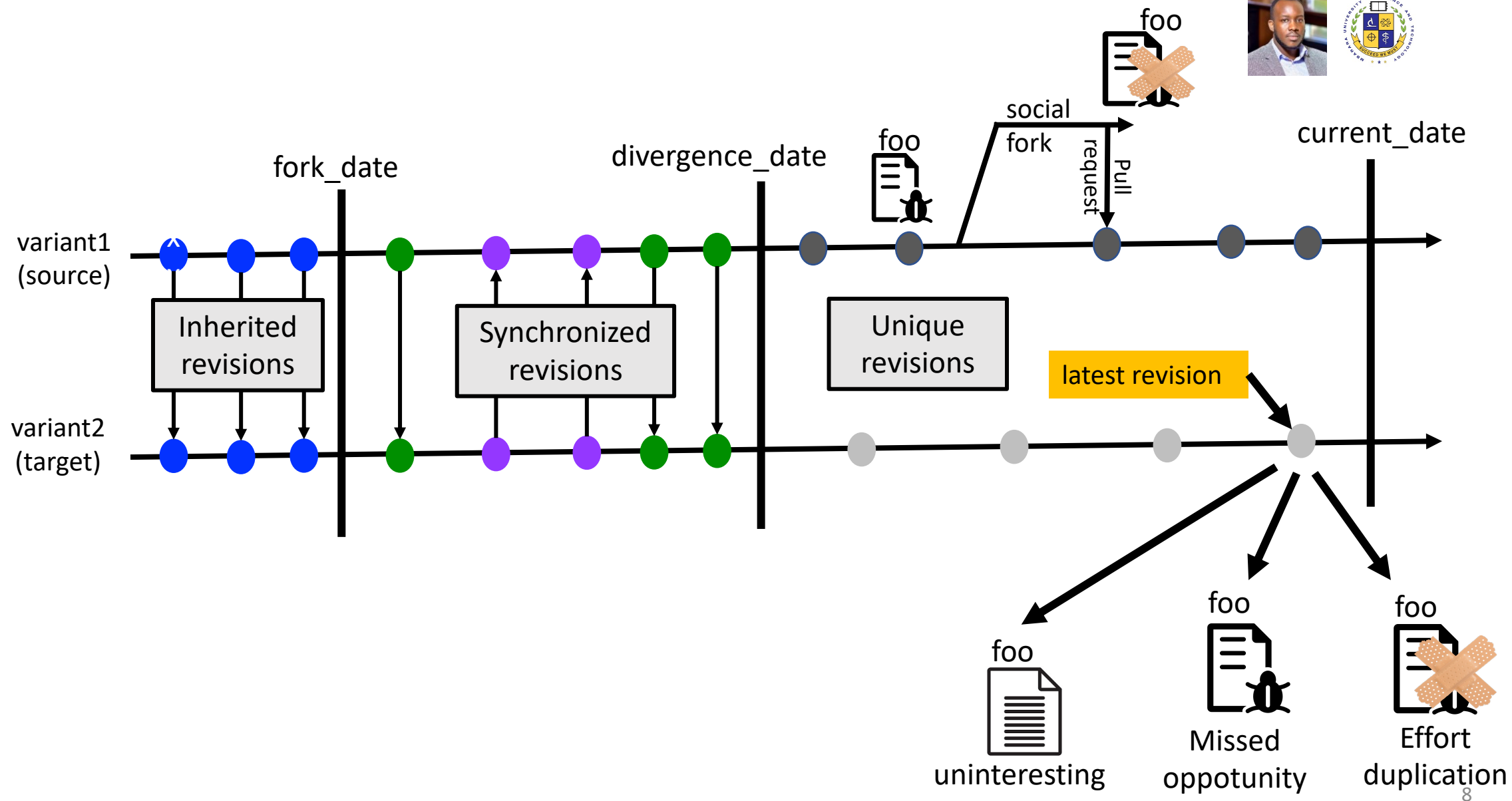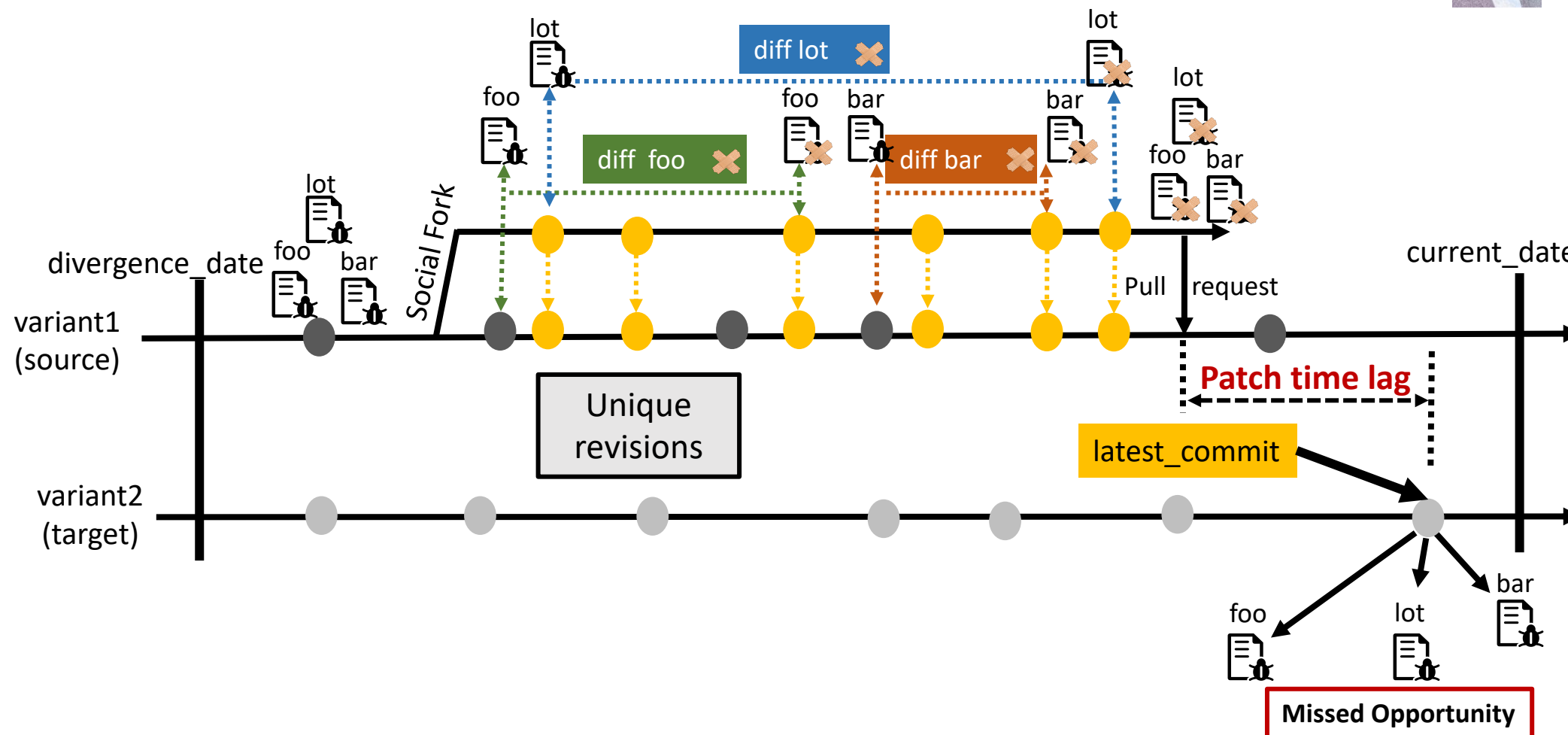
# Illustration of diverged variants – Zoomed patch

# Concrete Example – Missed Opportunity

## Buggy snippet from upstream

```
1              return;
2          }
3      } while (p < (uint16_t *)SYMVAL(__eeprom_workarea_end__));
4      flashend = (uint32_t)((uint16_t *)SYMVAL(__eeprom_workarea_end__) − 1);
5  }
```

← Buggy line

## Patched snippet from upstream

```
1              return;
2          }
3      } while (p < (uint16_t *)SYMVAL(__eeprom_workarea_end__));
4      flashend = (uint32_t)(p − 1);
5  }
```

← Patched line

## Diff for patch in upstream

```
1  @@ −363,7 +363,7 @@
2
3      } while (p < (uint16_t *)SYMVAL(__eeprom_workarea_end__));
4  −    flashend = (uint32_t)((uint16_t *)SYMVAL(__eeprom_workarea_end__) − 1);
5  +    flashend = (uint32_t)(p − 1);
```

} Hunk

## Latest_commit snippet from divergent fork

```
1              return;
2          }
3      } while (p < (uint16_t *)SYMVAL(__eeprom_workarea_end__));
4      flashend = (uint32_t)((uint16_t *)SYMVAL(__eeprom_workarea_end__) − 1);
5  }
```

← Buggy line

# Concrete Example – Effort Duplication

Buggy snippet from upstream

```
1        # http://ss64.com/nt/syntax-esc.html
2        _escape_re = re.compile(r'(?<!\^)[&<>]|(?<!\^)\^(?![&<>\^])')      ← Buggy line
3        _escaper = partial(_escape_re.sub, lambda m: '^' + m.group(0))
```

Patched snippet from upstream

```
1        # http://ss64.com/nt/syntax-esc.html
2        _escape_re = re.compile(r'(?<!\^)[&<>]|(?<!\^)\^(?![&<>\^])|(\|)')   ← Patched line
3        _escaper = partial(_escape_re.sub, lambda m: '^' + m.group(0))
```

Diff for patch in upstream

```
1        @@ -24,7 +24,7 @@
2
3        # http://ss64.com/nt/syntax-esc.html
4    -   _escape_re = re.compile(r'(?<!\^)[&<>]|(?<!\^)\^(?![&<>\^])')         ⎫
5    +   _escape_re = re.compile(r'(?<!\^)[&<>]|(?<!\^)\^(?![&<>\^])|(\|)')    ⎬ Hunk
6        _escaper = partial(_escape_re.sub, lambda m: '^' + m.group(0))       ⎭
```
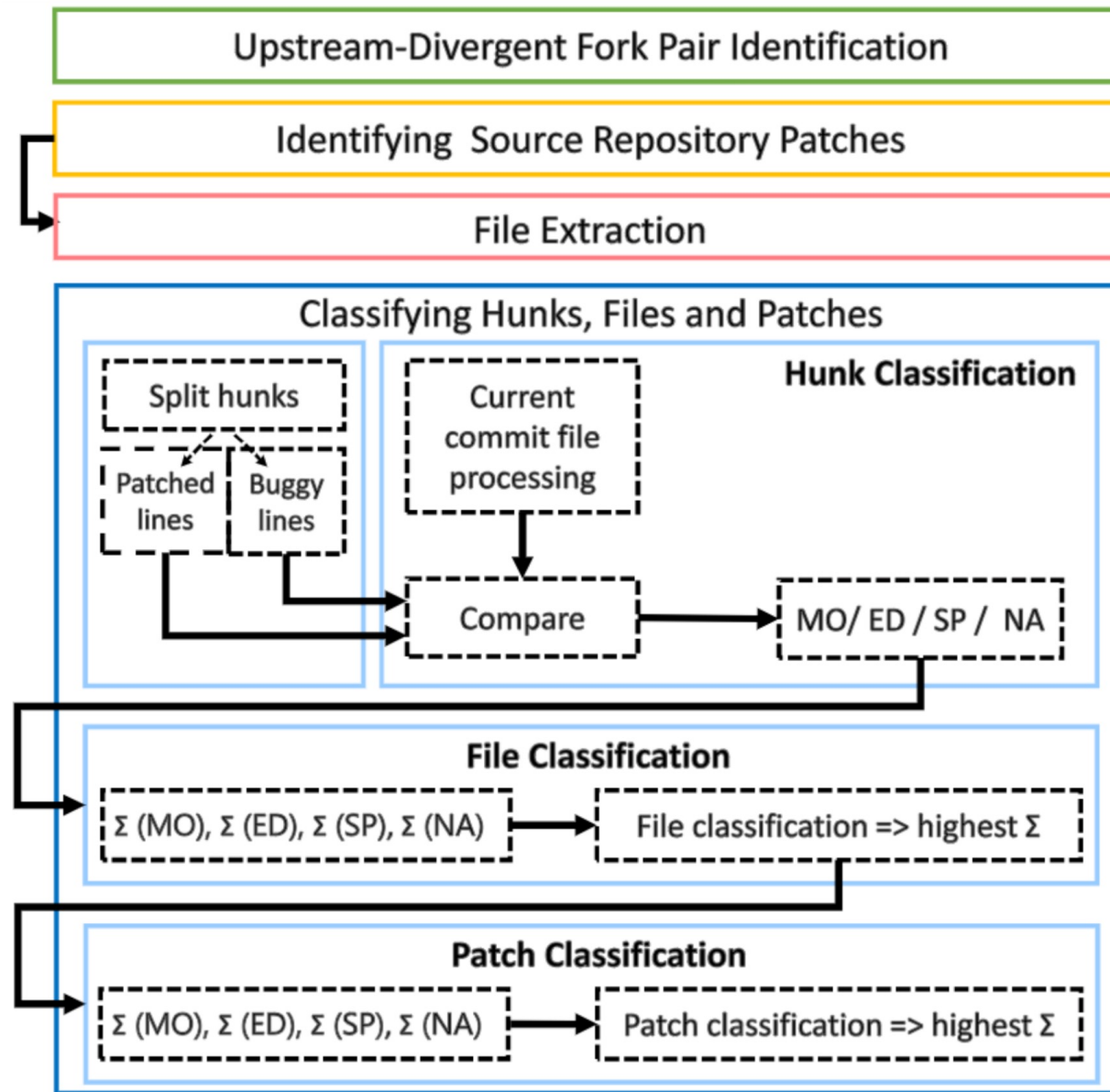
Latest_commit snippet from divergent fork

```
1        # http://ss64.com/nt/syntax-esc.html
2        _escape_re = re.compile(r'(?<!\^)[&<>]|(?<!\^)\^(?![&<>\^])|(\|)')   ← Patched line
3        _escaper = partial(_escape_re.sub, lambda m: '^' + m.group(0))
```

# Tool - PeReco

# Let us play with the tool