# USABLE PROGRAMMING TOOLS FOR EXPERIMENTAL BIOLOGISTS

**Justin Lubin**

Advisor: Sarah E. Chasins

EPIC Retreat, Spring 2022

Vaccines

Precision health

Genomic editing

Wet lab ⟷ Dry lab

# COMPONENT-BASED REFACTORING

# WET LAB LANGUAGE
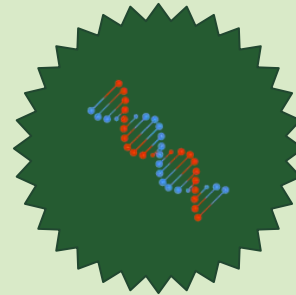
# COMPONENT-BASED REFACTORING

WET LAB LANGUAGE

# COMPONENT-BASED REFACTORING

1. Motivating example

2. Naïve approach
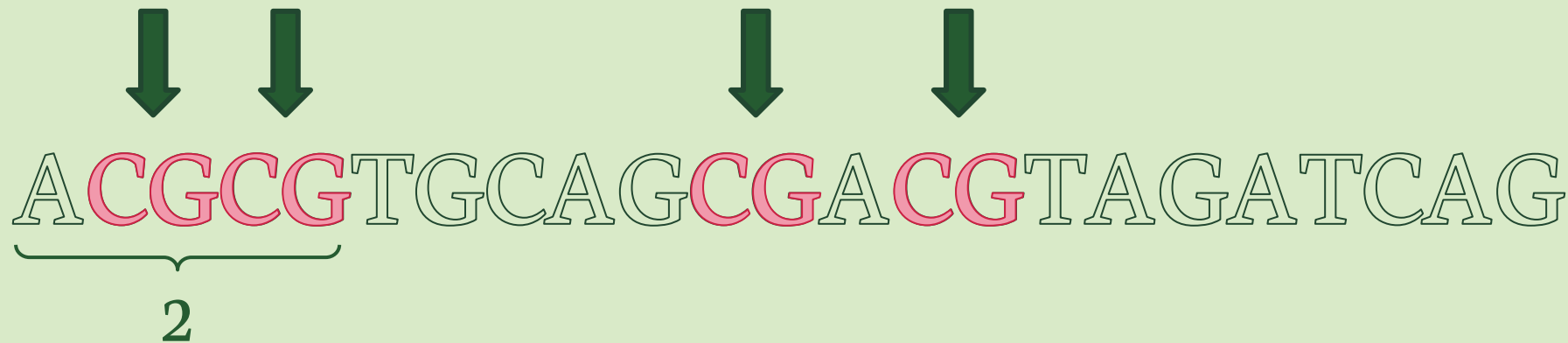
3. Our approach

# COMPONENT-BASED REFACTORING

ACGCGTGCAGCGACGTAGATCAG

ACGCGTGCAGCGACGTAGATCAG

2 2 1

ACGCGTGCAGCGACGTAGATCAG

2 2 1 1 0

# Straightforward solution:

Two nested for loops

*...but takes **>1 day** to run...*

# Faster solution:

```
np.concatenate([
  np.convolve(
    (seq == "C")[:-1] & (seq == "G")[1:],
    np.ones(window_size),
    "valid"),
  [0]])
```

## Straightforward solution:

Two nested for loops

*...but takes **>1 day** to run...*

## Faster solution:

```python
np.concatenate([
  np.convolve(
    (seq == "C")[:-1] & (seq == "G")[1:],
    np.ones(window_size),
    "valid"),
  [0]])
```

*...but takes **<15 min** to run!*

# COMPONENT-BASED REFACTORING

# COMPONENT-BASED REFACTORING

# COMPONENT-BASED REFACTORING

# COMPONENT-BASED REFACTORING

1. Motivating example

2. Naïve approach

3. **<u>Our approach</u>**

```
len(x)
x[i] * y[i]
```

```
len(mul(x, y))
mul(x, y)[i]
```

# CANONICALIZATION

`len(mul(x, y))` ➡️ `len(x)`

`mul(x, y)[i]` ➡️ `x[i] * y[i]`
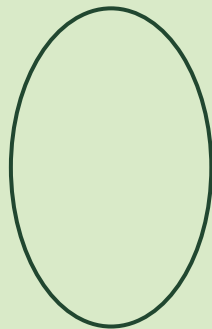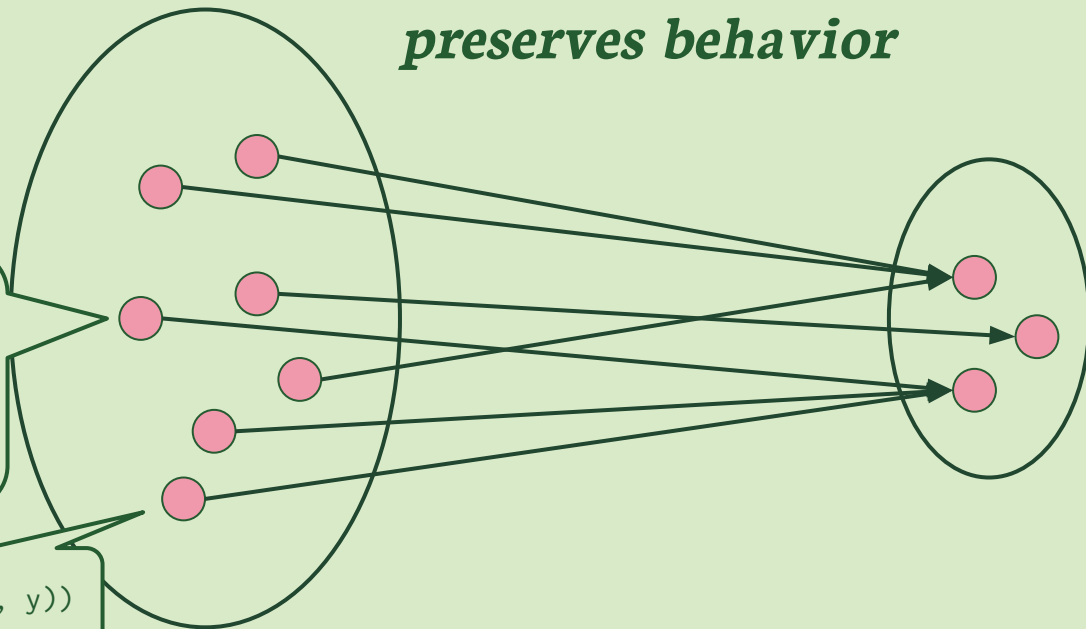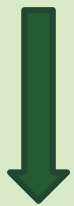
*Canonicalize*

*Canonicalize*

*preserves behavior*

```
def dot(x, y):
    total = 0
    for i in range(len(x)):
        total += x[i] * y[i]
    return total
```
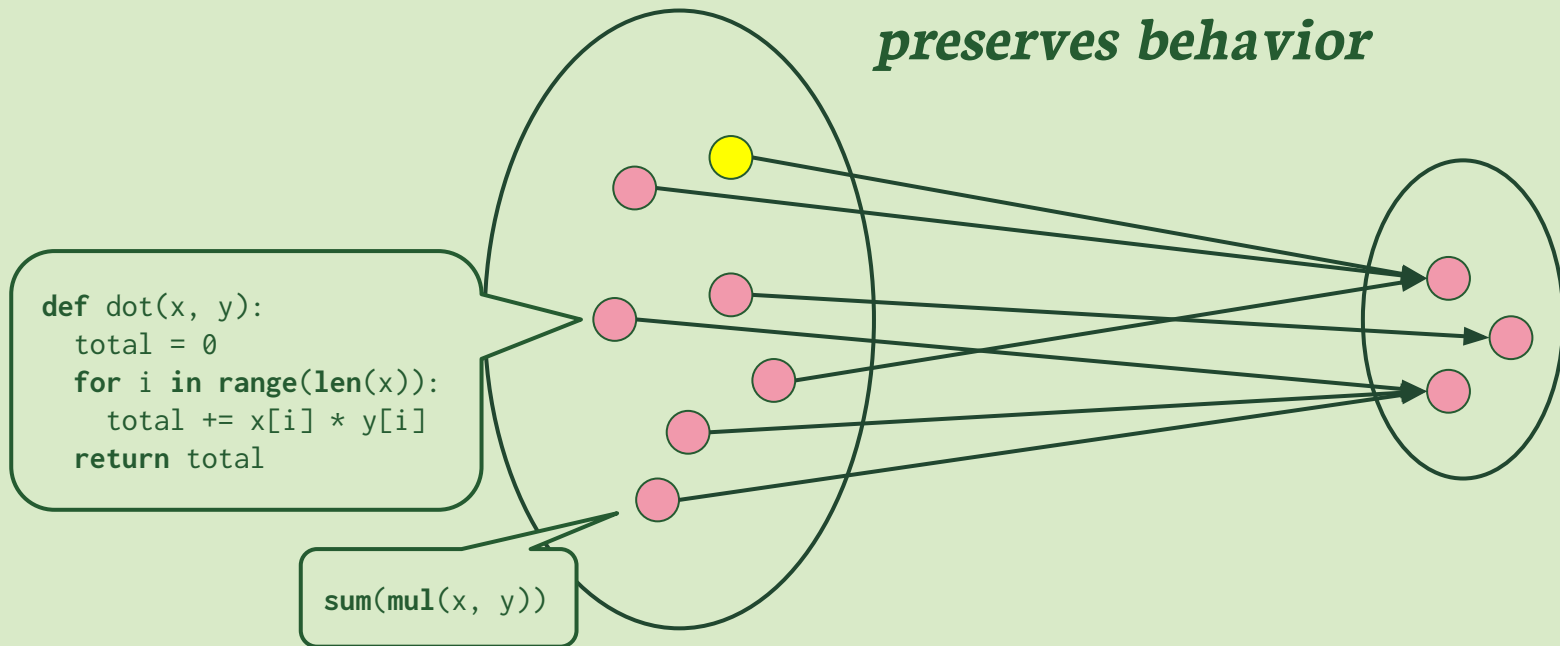
```
sum(mul(x, y))
```

*Canonicalize*

*preserves behavior*

```
def dot(x, y):
    total = 0
    for i in range(len(x)):
        total += x[i] * y[i]
    return total
```
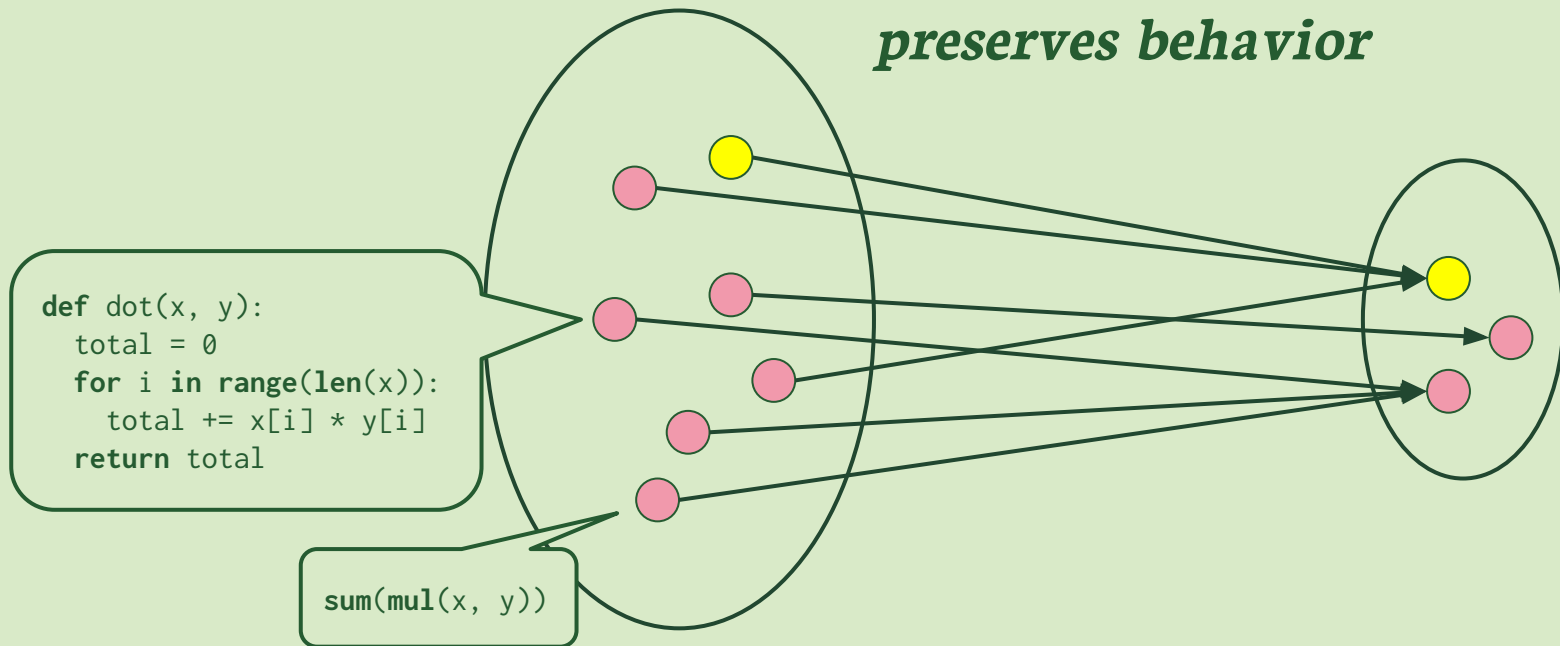
```
sum(mul(x, y))
```
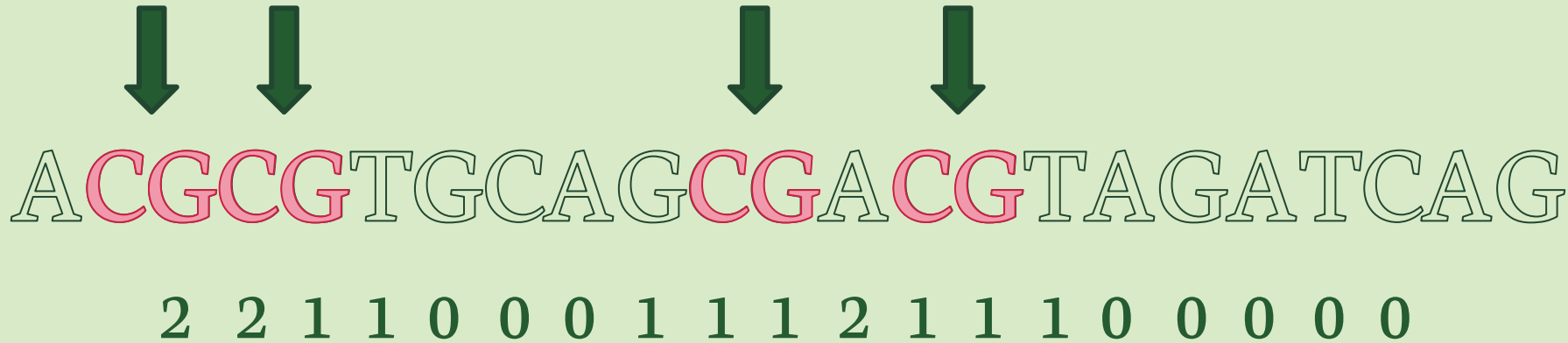
For many programs,

# KEY IDEA

# KEY IDEA

For many programs,

we can check
**extensional equality**

— *by* —

**syntactic equality**
modulo
**canonicalization**

A C G C G T G C A G C G A C G T A G A T C A G

2 2 1 1 0 0 0 1 1 1 2 1 1 1 0 0 0 0 0

## Straightforward solution:

Two nested for loops

*...takes **>1 day** to run...*

## Faster solution:

```
np.concatenate([
  np.convolve(
    (seq == "C")[:-1] & (seq == "G")[1:],
    np.ones(window_size),
    "valid"),
  [0]])
```

*...takes **<15 min** to run!*

**Jeremy Ferguson**                    **Kevin Ye**                    **Jacob Yim**