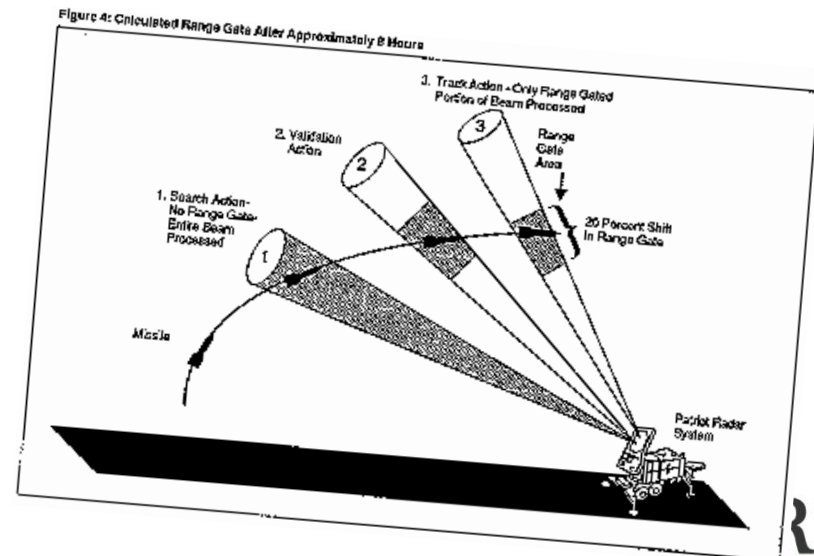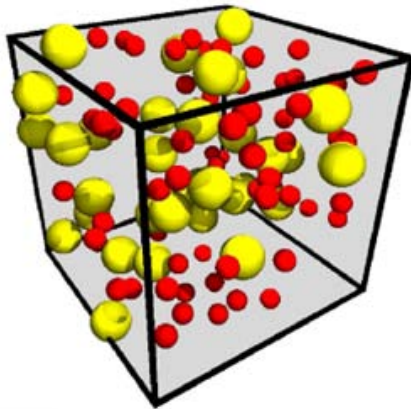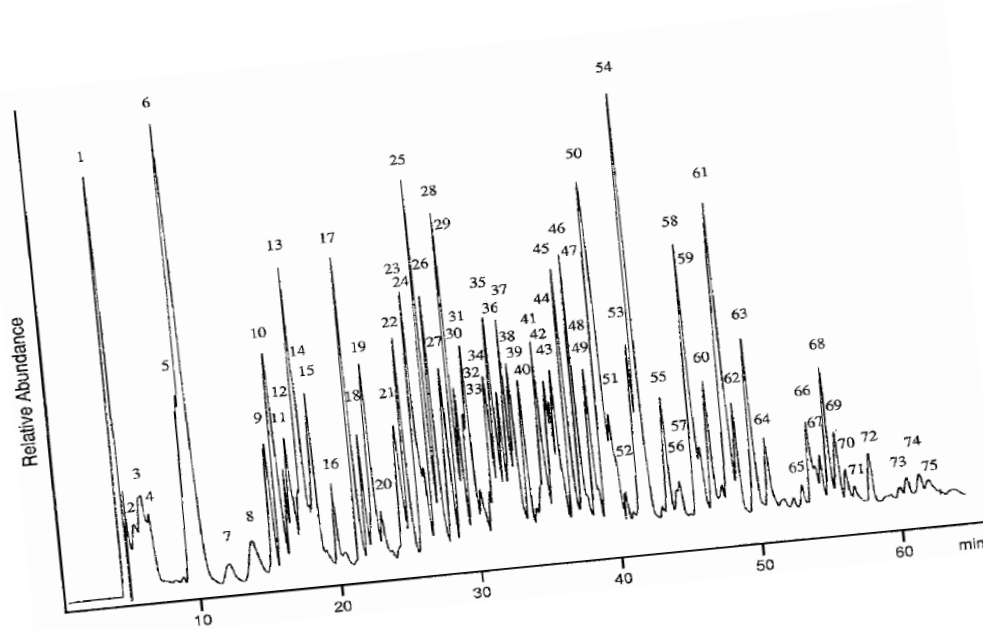# Reliable Data Processing Enabled by Program Analysis

Xiangyu Zhang

December 22, 2013 @ ISHCS, PKU

PURDUE
UNIVERSITY

# Data processing is becoming increasing important.

# Errors in Data Processing

*external errors*

*internal errors*

```
    _kbucket_index(self, key)
    """
    Returns the index of the k
    string.
    """
    if isinstance(key, str):
        key = hex_to_long(key)
    # Bound check for key too
    if key < 0:
        raise ValueError('Key
    i = 0
    for bucket in self._bucke
        if bucket.key_in_rang
            return i
        else:
            i += 1
    # Key was too big given
    raise ValueError('Key ou
```

PURDUE
UNIVERSITY

# External Errors

```
1    float x, z;
2    x = input();
3    z = f(x);

4    if (z > 0.5)
5        printf ("hit");
6    else
7        printf ("miss");
```

$X = 100.0;$

1. How would the program output change if input $x$ is uncertain?

e.g. $x \in [50.0, 150.0]$

# Internal Errors

```
1   float x, z;
2   x = input();
3   z = f(x);

4   if (z > 0.5)
5       printf ("hit");
6   else
7       printf ("miss");
```

X = 100.0;

2. Are the computed results reliable?

miss

PURDUE
UNIVERSITY

# Errors in Data Processing

- External errors
  - Also known as data uncertainty problem
  - Existing techniques
    - query-based uncertain data processing e.g. [R. Jampani, SIGMOD 2008], [S. Singh, ICDE 2008] and etc.
    - Interval analysis
    - Automatic differentiation

- Internal errors
  - Existing work include interval analysis, using high precision

- Existing solutions are hardly applicable or usable -- too expensive, too many false positives

- Errors may get propagated and magnified, leading to unreliable output.
  - We call it the *instability* problem.

PURDUE
UNIVERSITY

# Outline

- Overview
- External Errors
  - White-box sampling (OOPSLA 2012)
- Internal Errors
  - On-the-fly detection of instability problems (OOPSLA 2013)

**PURDUE**
UNIVERSITY

# White Box Sampling --External Errors

```
1  float x, z;                    X = 100.0;
2  x = input();
3  z = f(x);

4  if (z > 0.5)
5      printf ("hit");
6  else
7      printf ("miss");
```
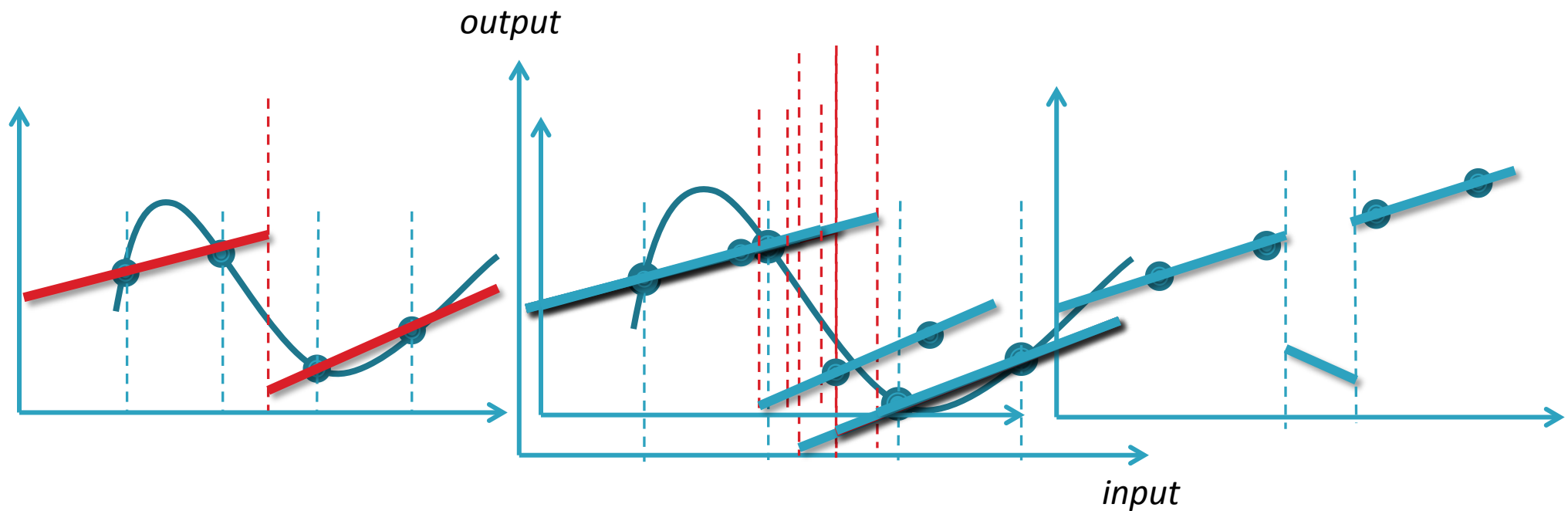
> 1. How would the program output change if input $x$ is uncertain?
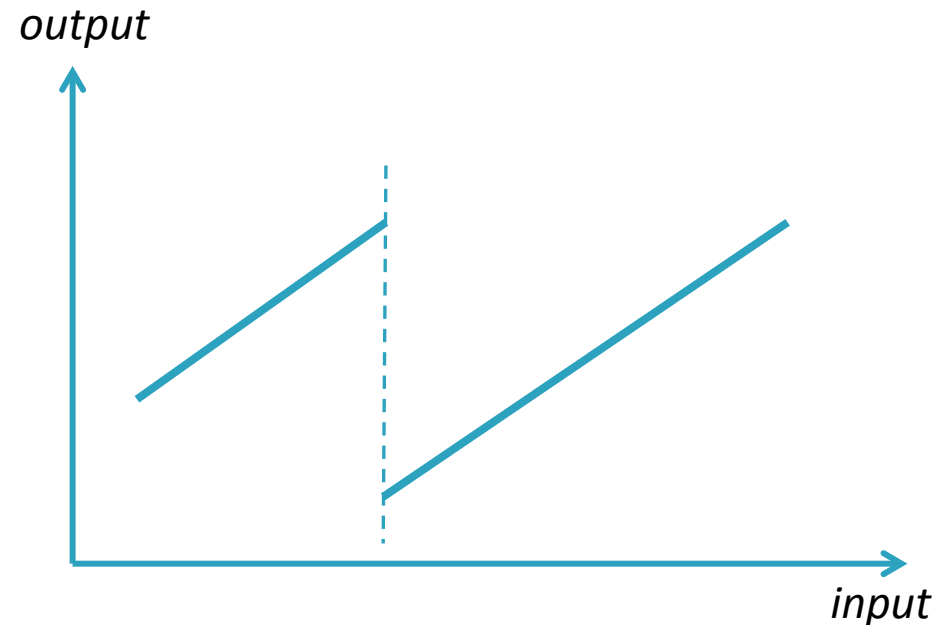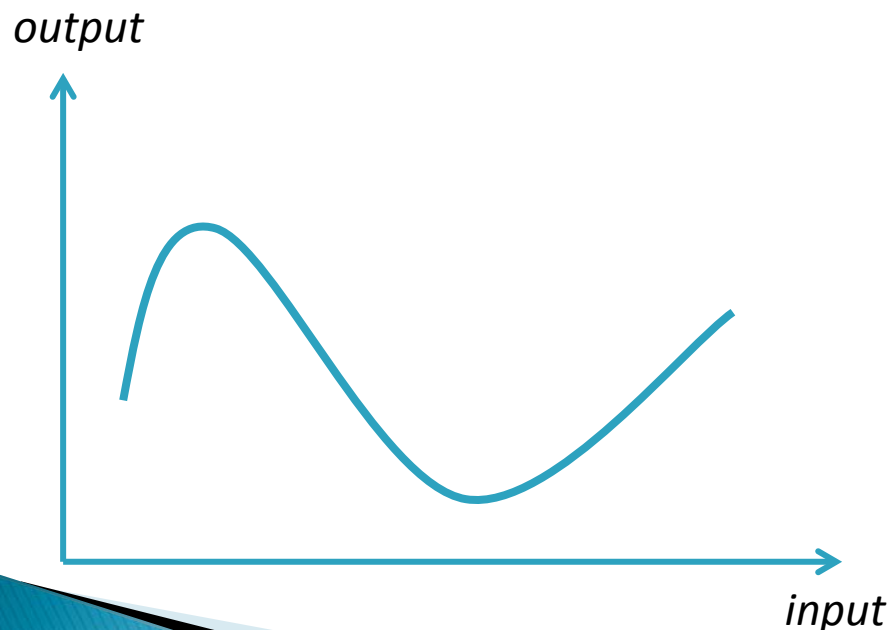>
> e.g. $x \in [50.0, 150.0]$

PURDUE
UNIVERSITY

# Monte Carlo (MC) methods

- Sampling-based *Monte Carlo method* is effective, yet imprecise.

# Our Idea

- MC sampling guided by program analysis (few samples and sampling at the critical places)
  - Use dynamic analysis to predict output continuity
  - Perform demand driven sampling based on continuity

# White Box Sampling - Intuition

# A Running Example

```
1   x = sample(1.5);
2   y = (int) x;
3   if (x < 1.0)
4     o = 1 + y;
5   else
6     if (t(x) > 0.3)
7       o = 0.3;
8     else
9       o = 0.75;
```