

# Precise Dynamic Data Flow Tracking with Proximal Gradients

Gabriel Ryan<sup>1</sup>, Abhishek Shah<sup>1</sup>, Dongdong She<sup>1</sup>, Koustubha Bhat<sup>2</sup>, Suman Jana<sup>1</sup>

<sup>1</sup>Columbia University

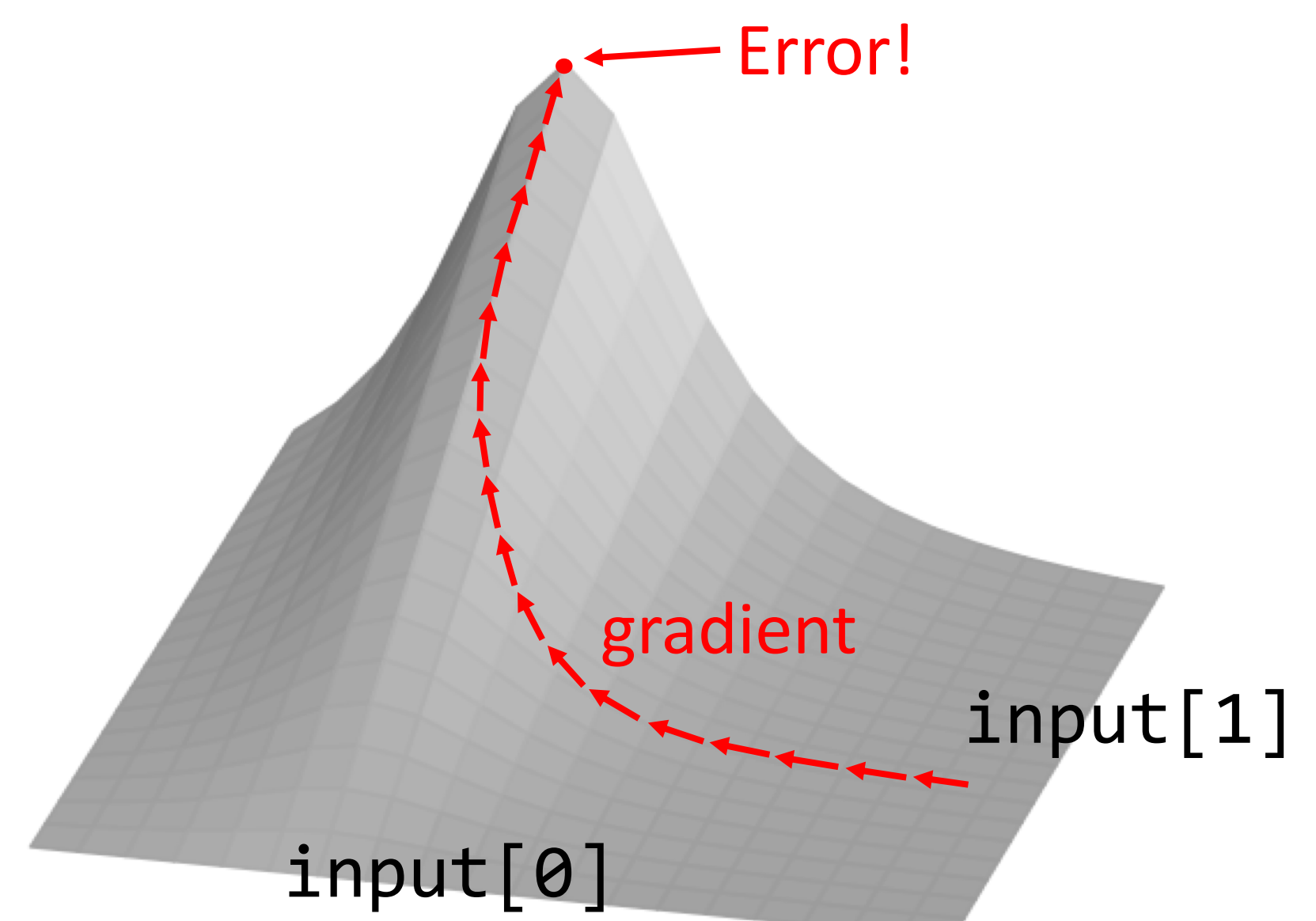
<sup>2</sup>Vrije Universiteit

gabe@cs.columbia.edu, as5258@columbia.edu, ds3619@columbia.edu, k.bhat@vu.nl, suman@cs.columbia.edu



```
taint source: x
taint sink: y
// x is a 4 byte int
int x = 0x12345678;
for (int i=0; i<6; i++){
    y[i] = x;
    x = x<<8;
}
```

y					
1	1	1	1	1	1
Taint to y from x					
1	256	2 <sup>16</sup>	2 <sup>24</sup>	0	0
Gradient of y wrt. x					



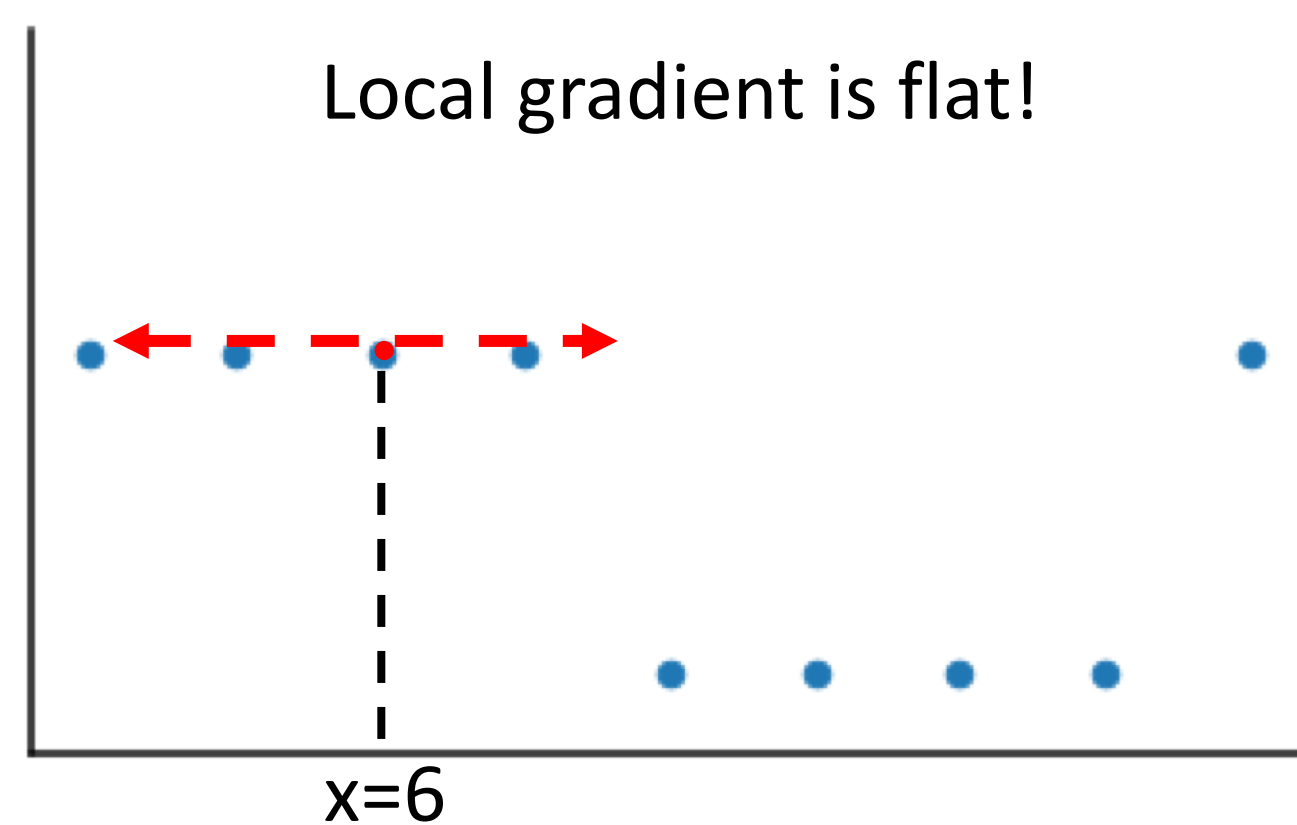
Gradient can be used to guide searches over large input spaces.

Figure 1. Gradient avoids approximation errors and identifies where  $x$  has the greatest effect on  $y$ .

Gradient provides more precise and meaningful information about program behavior than taint based dataflow tracking.

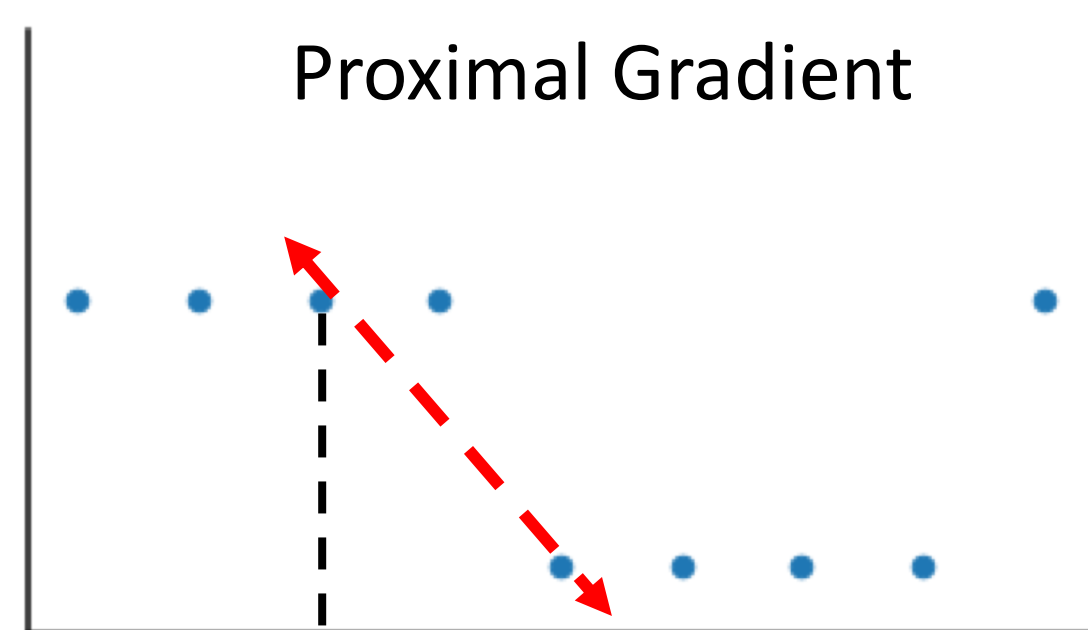
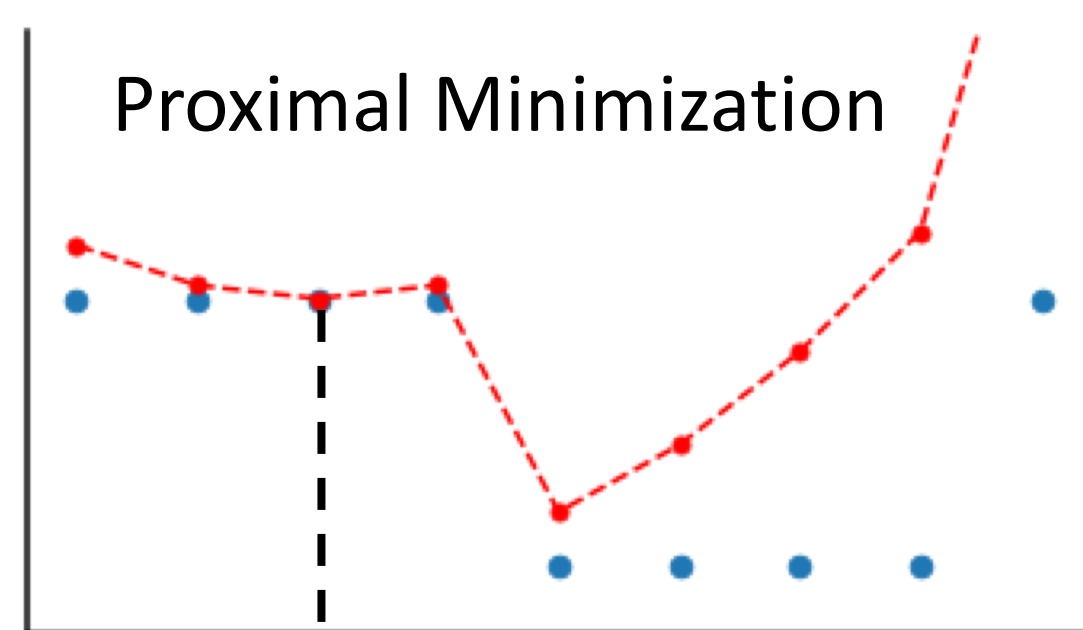
## Problem: Programs are not Differentiable

```
int f(int x) {
    return x & 4;
}
```



Programs contain nonsmooth operations like Bitwise And, making it impossible to compute gradients normally.

## Solution: Nonsmooth Optimization



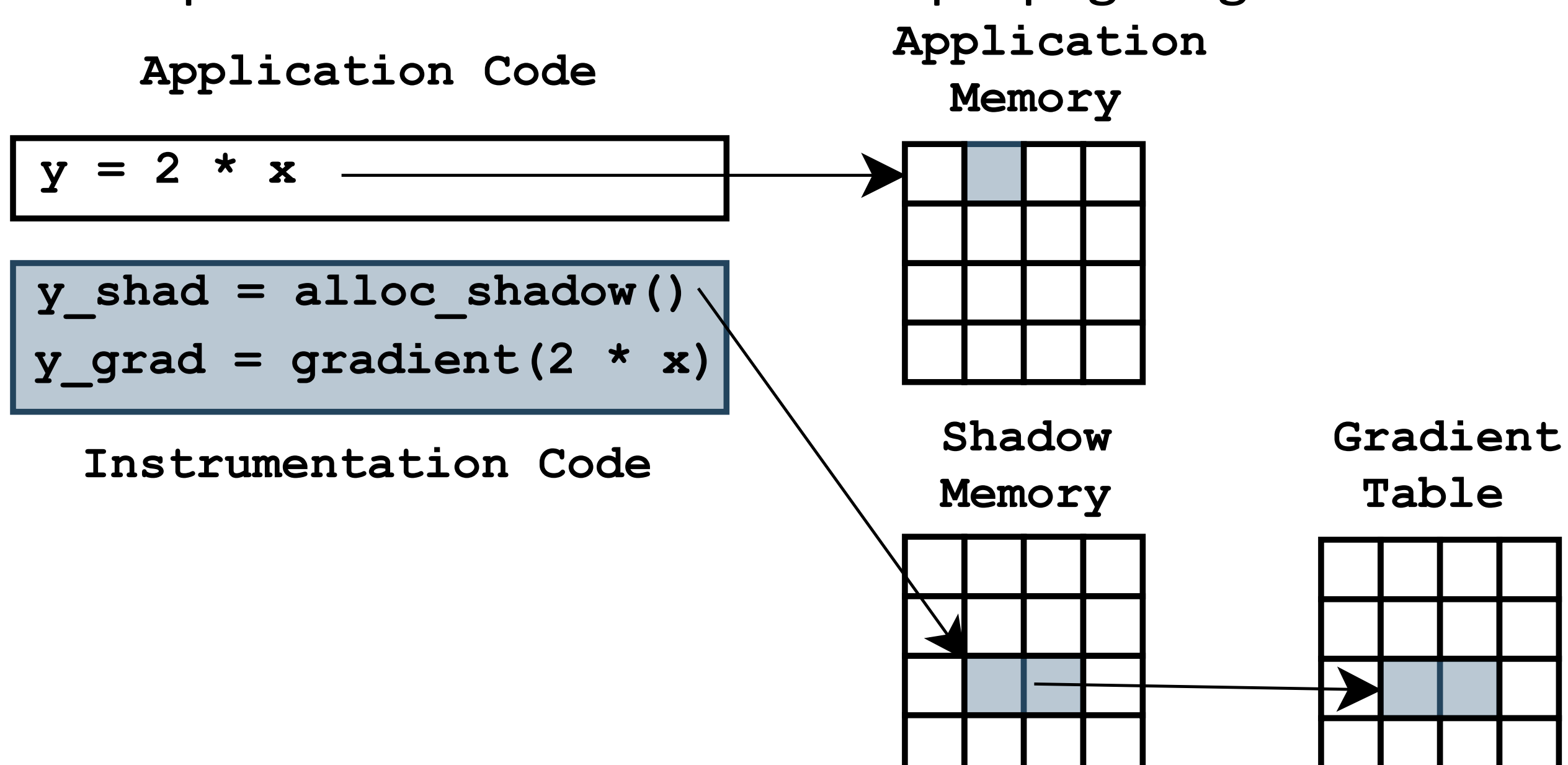
Proximal Operator finds nearby minimum

Proximal Gradient is computed from nearby minimum

Proximal Gradients are an application of nonsmooth optimization methods that base the gradient on the local minimum of a function.<sup>3</sup>

## Implementation:

- Implemented as LLVM Sanitizer, based on DataFlowSanitizer
- Each address has associated shadow memory with label
- Each operation is instrumented to propagate gradient

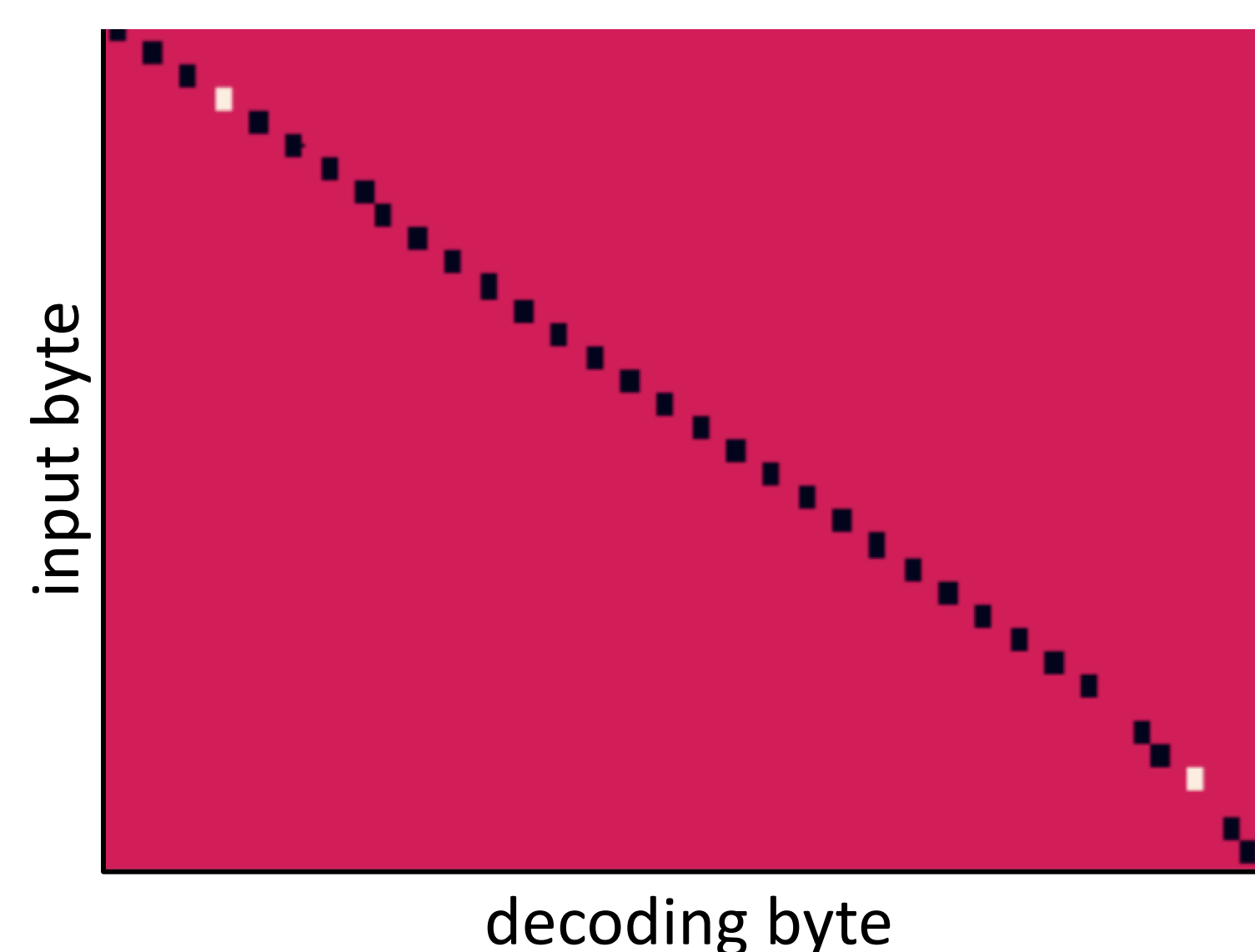
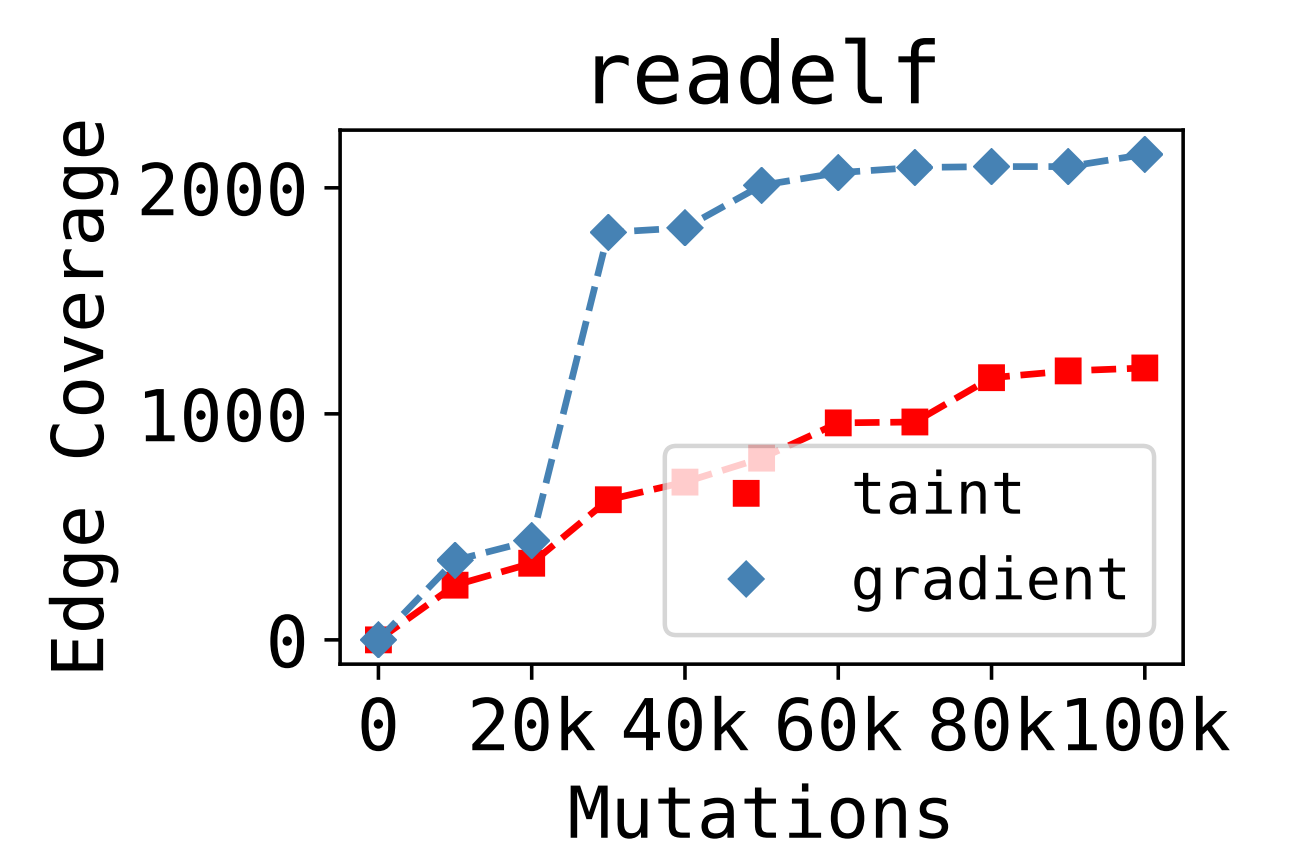
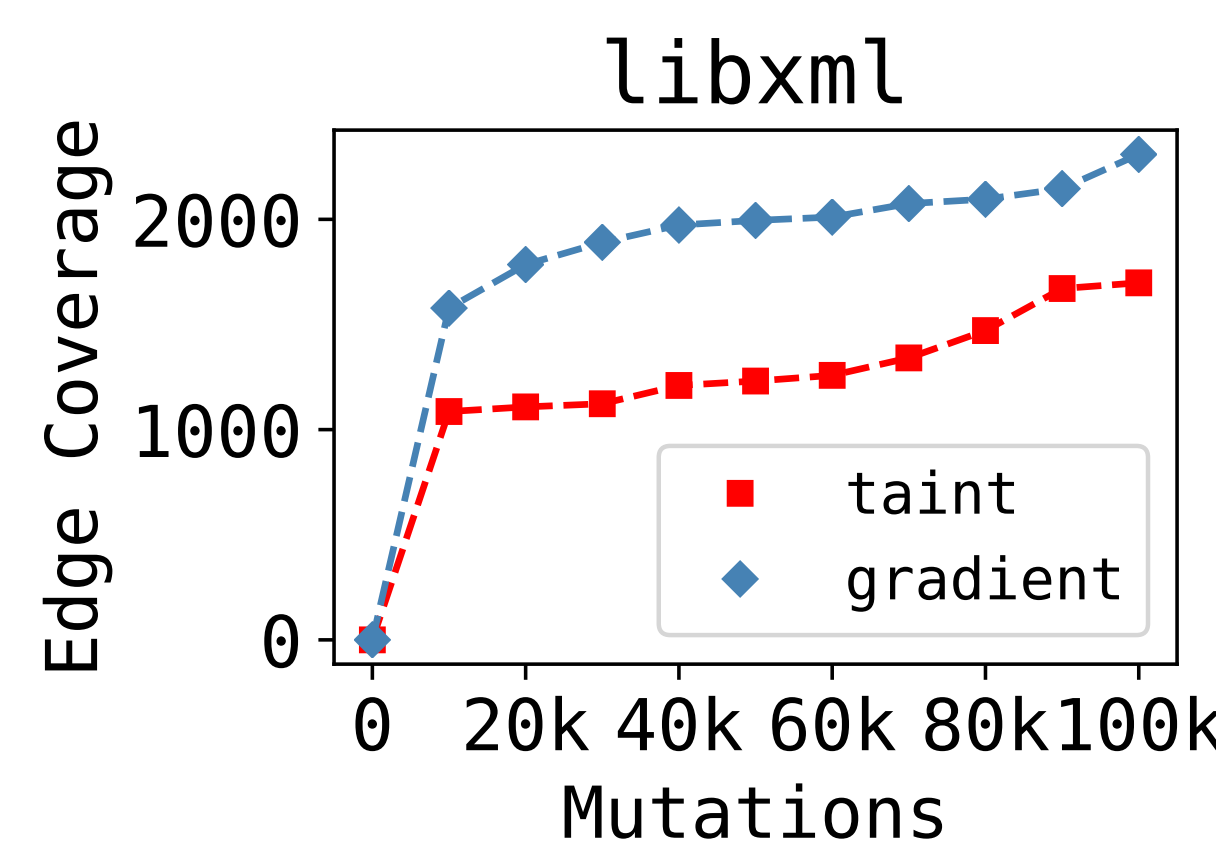


## Evaluation:

- Achieved up to 39% better dataflow precision than state of the art method (Taint Analysis) on 7 programs.

Program	Taint			Gradient		
	Precision	Recall	F1	Precision	Recall	F1
minigzip	0.55	0.86	0.68	0.94	0.71	<b>0.81</b>
djpeg	0.63	0.73	0.68	0.96	0.61	<b>0.74</b>
mutool	1.0	0.01	<b>0.02</b>	1.0	0.01	<b>0.02</b>
xmllint	0.97	0.39	<b>0.56</b>	0.97	0.39	<b>0.56</b>
readelf	0.17	0.95	0.28	0.18	0.93	<b>0.30</b>
objdump	0.77	0.80	0.78	0.94	0.79	<b>0.85</b>
strip	0.60	0.83	0.70	0.88	0.79	<b>0.84</b>

- Overhead is comparable to Taint Analysis on average (<5%).
- Example application: Dataflow guided fuzzing.
  - Increase of 57% edge coverage on 2 programs



Gradient Analysis identifies which input bytes are used to look up each Huffman decoding index in gzip decompression.

## References:

- [1] LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. "Deep learning." *nature* 521.7553 (2015): 436.
- [2] Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." *Proceedings of the IEEE International Conference on Computer Vision*. 2017.
- [3] Parikh, Neal, and Stephen Boyd. "Proximal algorithms." *Foundations and Trends® in Optimization* 1.3 (2014): 127-239.