Microsoft Research
*Software Analytics* Group

# *Context-Sensitive Delta Inference*
# *for Identifying Workload-Dependent*
# *Performance Bottlenecks*

Xusheng Xiao[1], Shi Han[2], Dongmei Zhang[2], **Tao Xie**[1,3]

[1]North Carolina State University
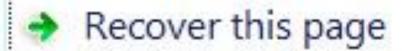[2]Microsoft Research Asia
[3]University of Illinois at Urbana-Champaign
[1]xxiao2@ncsu.edu, [2]{shihan, dongmeiz}@microsoft.com, [1,3]taoxie@illinois.edu

# Windows Internet Explorer

## A webpage is not responding on the following website: facebook.com

You can wait for the webpage to respond, or choose one of the following options:

➡ Recover this page

➡ Close this page

⌄ See details

# End Program - ▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓ - Win...

This program is not responding.

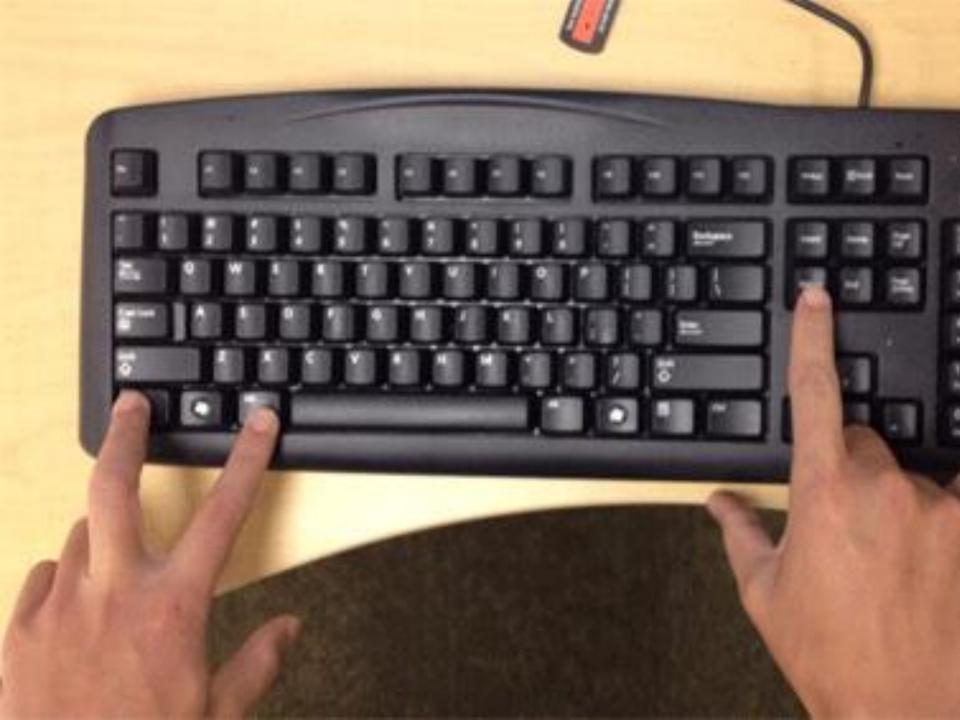To return to Windows and check the status of the program, click Cancel.

If you choose to end the program immediately, you will lose any unsaved data. To end the program now, click End Now.
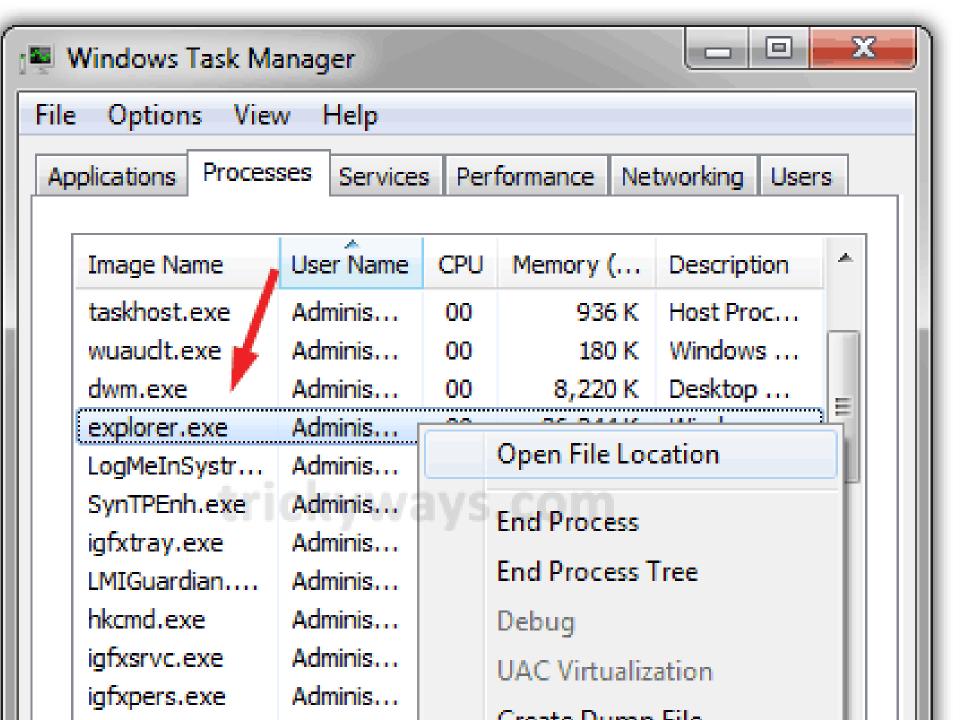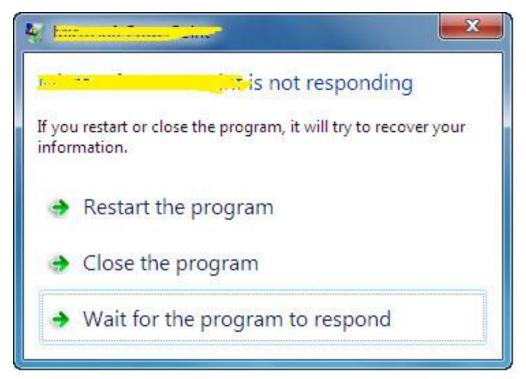
[ End Now ]     [ Cancel ]

# Windows Task Manager

File  Options  View  Help

| Applications | Processes | Services | Performance | Networking | Users |

| Image Name | User Name | CPU | Memory (... | Description |
|------------|-----------|-----|-------------|-------------|
| taskhost.exe | Adminis... | 00 | 936 K | Host Proc... |
| wuaudt.exe | Adminis... | 00 | 180 K | Windows ... |
| dwm.exe | Adminis... | 00 | 8,220 K | Desktop ... |
| explorer.exe | Adminis... | | | |
| LogMeInSystr... | Adminis... | | | |
| SynTPEnh.exe | Adminis... | | | |
| igfxtray.exe | Adminis... | | | |
| LMIGuardian.... | Adminis... | | | |
| hkcmd.exe | Adminis... | | | |
| igfxsrvc.exe | Adminis... | | | |
| igfxpers.exe | Adminis... | | | |

Open File Location

End Process

End Process Tree

Debug

UAC Virtualization

Create Dump File

trickyways.com

# Performance Problems

- Widely exist in released software
  - Mozilla developers fixed **5–60** perf bugs **monthly** over the past 10 years [Jin et al. PLDI 12]



**Software Hangs** in daily tools: *file managers, office tools, browsers,*

*...*

# Software Hangs

- Three major categories
  - Correctness, e.g., infinite loops
  - Blocking operations, e.g., sending files
  - Expensive operations

contribute to **27%** of 233 hang bugs

[Song et al. DSN 2010]

*With **root** of disk ( C:\ ) **selected** in drop-down path selector, attempting to **enable Flat View** under the top-level View menu causes 7-Zip to hang …*

*The process had to be forcibly stopped using Windows Task Manager.*

7-Zip File Manager

# Workload-Dependent Performance Bottlenecks (WDPB)

- Expensive operations depending on input workloads, e.g., data processing
- **Insight**: caused often by **workload-dependent loops** with **expensive** operations, e.g.,
  - Temp-object creation/destruction
  - File I/O
  - UI updates
- Fixed by
  - Spawning new threads
  - Limiting workload/processing sizes

# Target Problem:
# WDPB-Loop Prediction

Traditional
**Single**-Profile Setting
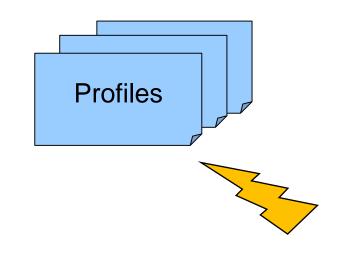
Target
**Multi**-Profile Setting

Profile

Profiles
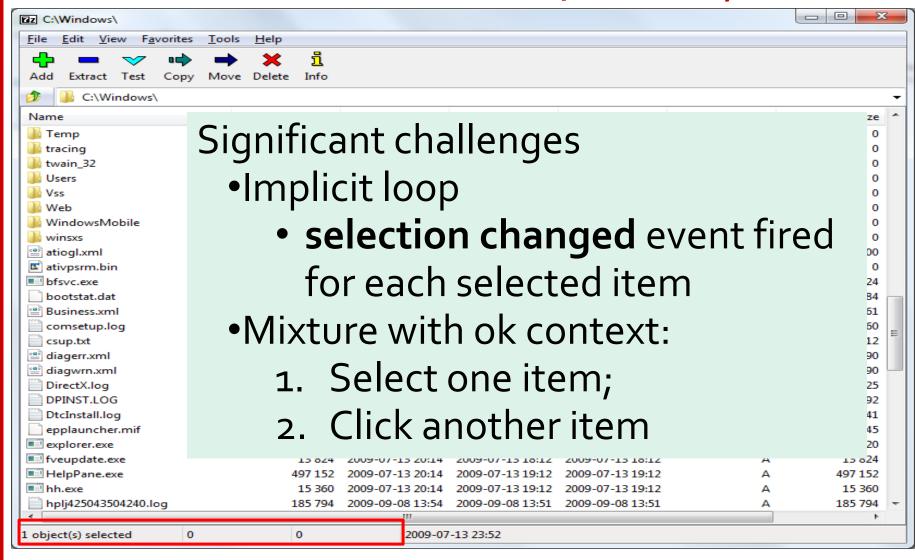
**Test Generation:**
Hard to know
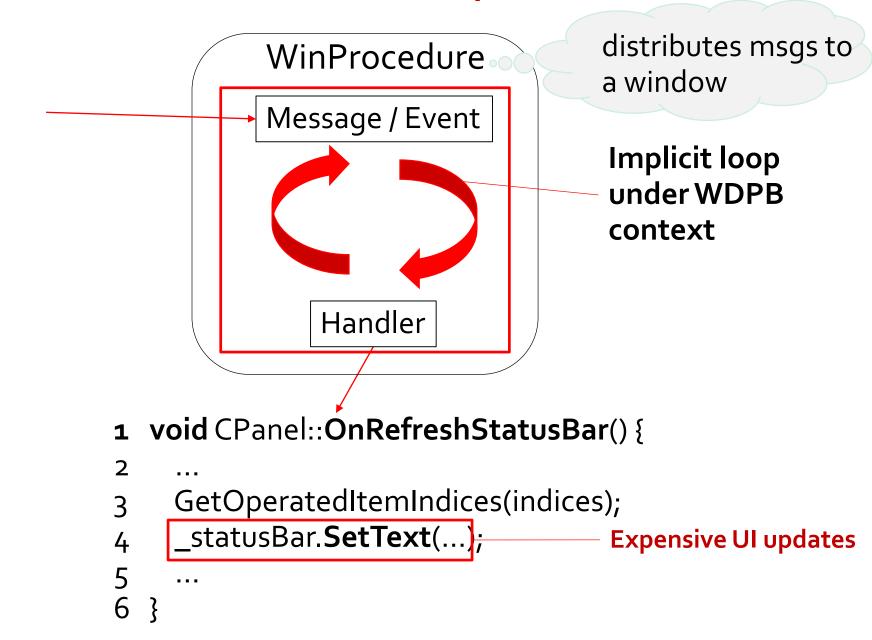triggering workload

**Test Oracle:**
Hard to know how
slow is slow enough

# Example WDPB: 7-Zip File Manager

## WDPB Context: 1. Select all items; 2. Click any item



Significant challenges
- Implicit loop
  - **selection changed** event fired for each selected item
- Mixture with ok context:
  1. Select one item;
  2. Click another item

# Inside the Example WDPB

WinProcedure

distributes msgs to a window

Message / Event

Implicit loop under WDPB context

Handler

```
1   void CPanel::OnRefreshStatusBar() {
2       …
3       GetOperatedItemIndices(indices);
4       _statusBar.SetText(…);
5       …
6   }
```

Expensive UI updates

# Background: *StackMine*

## Perf Debugging in the **Large**

Pattern Matching

Bug update

Problematic Pattern Repository

Bug Database

Network

Bug filing

How many issues are still unknown?

Key to issue discovery

Trace Storage

Bottleneck of scalability

Trace collection

Which trace file should I investigate first?

Trace analysis

# Proposed Approach: DeltaInfer
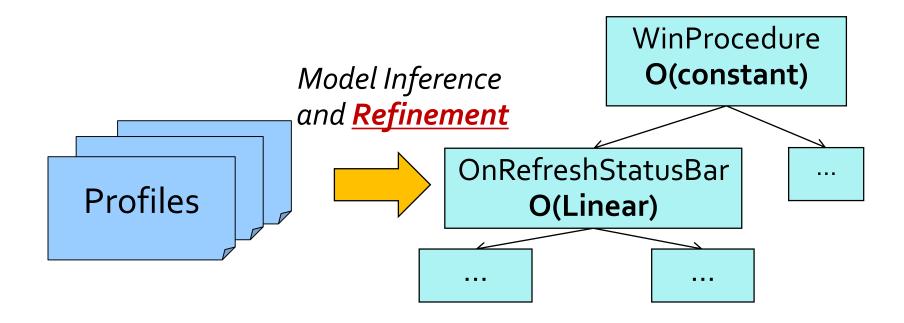## **Context-Sensitive** Delta Inference

- To gain the power of prediction

> ***Temporal*** *Inference: differences between* ***executions*** *as complexity models*

- To identify WDPB loops

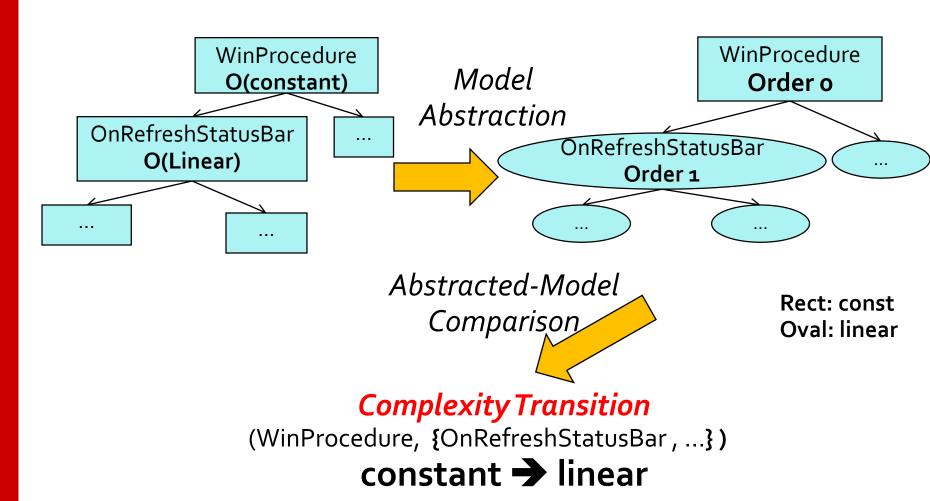> ***Spatial*** *Inference: differences between* ***program locations*** *as WDPB loops*

# Insight: Temporal Inference

*Model Inference and **Refinement***

Profiles → WinProcedure **O(constant)**

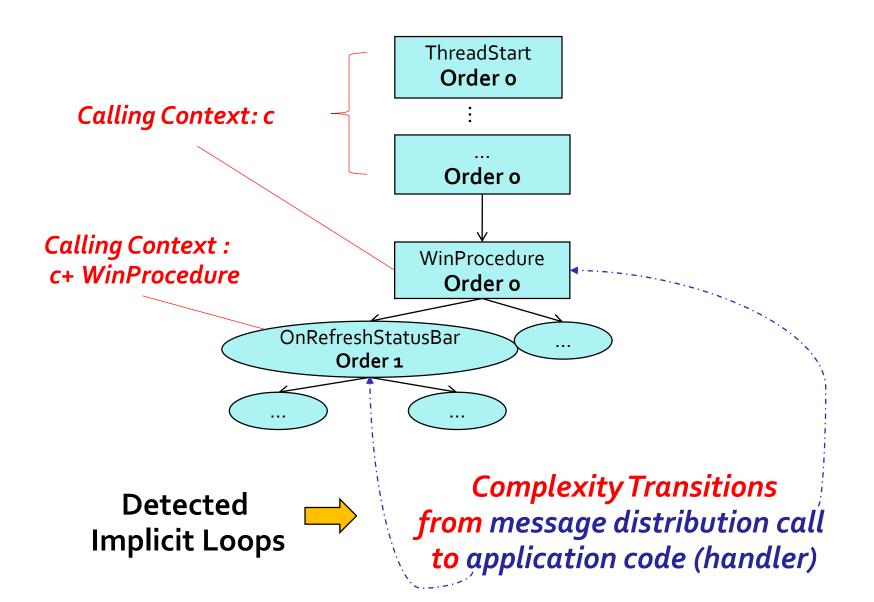OnRefreshStatusBar **O(Linear)**

...

...

...

Regression learning + **refinement** to produce **complexity models** of program locations

# Insight: Spatial Inference

WDPB loop **raises** complexity models of inside-loop locations to **higher order**

WinProcedure
**O(constant)**

OnRefreshStatusBar
**O(Linear)**

...

...

...

*Model Abstraction*

WinProcedure
**Order 0**

OnRefreshStatusBar
**Order 1**

...

...

...

*Abstracted-Model Comparison*

**Rect: const**
**Oval: linear**

*Complexity Transition*
(WinProcedure, {OnRefreshStatusBar , ...} )
**constant ➔ linear**

# Insight: Context Sensitivity

ThreadStart
**Order 0**

*Calling Context: c*

...
**Order 0**

WinProcedure
**Order 0**

*Calling Context :
c+ WinProcedure*

OnRefreshStatusBar
**Order 1**

...

...

...

**Detected
Implicit Loops**

*Complexity Transitions
from message distribution call
to application code (handler)*

# Overview of DeltaInfer

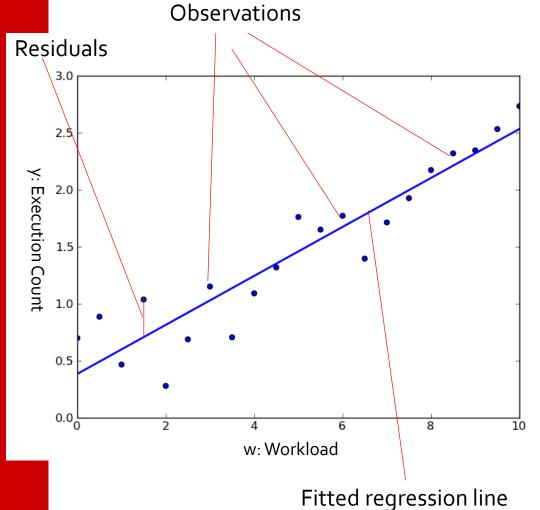# Workload Generation & Execution

Example scenario: *open a file in text editor*

- Performance metrics
  - execution time
- Performance-relevant workload parameters
  - **# of lines** (focused parameter)
    - Rep value range (RVR): [1, 1280]
    - Initial value/variations (sorted/random inputs)
  - # of character

*Example workloads*

  - # of lines (100, 200, ... , 500)
  - # of character (100 chars for a line)

# Model Inference



Observations

Residuals

y: Execution Count

w: Workload

Fitted regression line

- Linear Regression
  - $y = A + Bw$
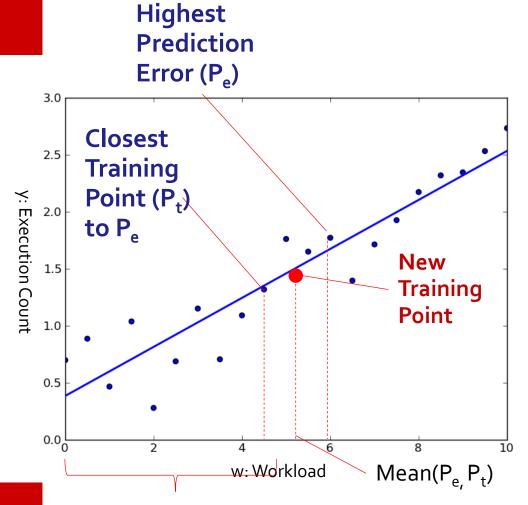- Power-law Regression
  - $y = Aw^B$

- Quality of the model
  - Correlation coefficient $R^2$

# Model Validation

- Model validation measures
  - relative prediction error of inferred models

- Example validation workloads:
  open a file in text editor
  - Validation value range (VVR): **[1,2560]**
  - Guideline: >= **2** times larger than the RVR
  - Caveat: too **large** RVR is **not** cost-effective

# Iterative Refinement



**Iterate** till

- Accuracy acceptable
- Improvement < threshold

**Select** a new workload

- Rationale: new workload at **highest-prediction-error areas** improves most

# Overview of DeltaInfer

**Temporal** Inference

**Spatial** Inference

Model Inference & Refinement

Model Abstraction

Models

Profiles          Models

Abstracted Models

**Initial Workloads**

Workload Generation & Execution

Abstracted Model Comparison

**Complexity Transitions**

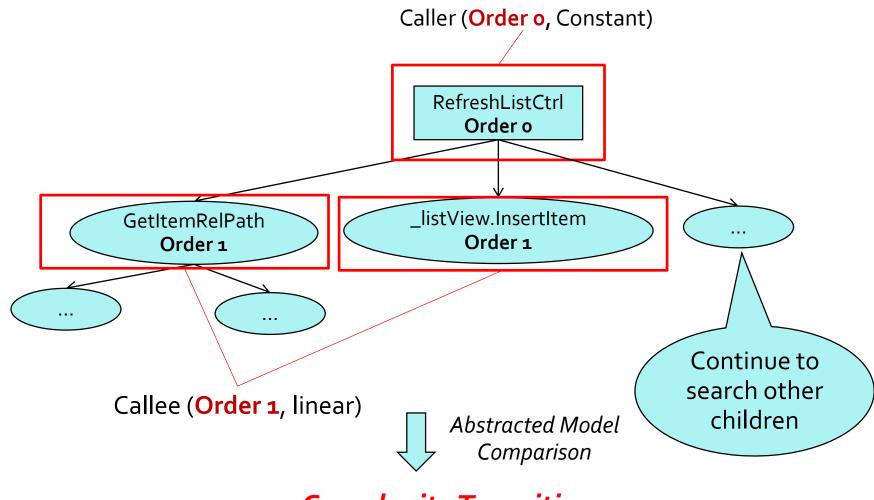# Model Abstraction: Complexity Orders



- Linear model ($y = A + Bw$)
  - **1**, if B > 0
  - **0**, otherwise

- Power-law model ($y = Aw^B$)
  - ***Round(B)***

- Model w/ $R^2$ below threshold$_{R2}$ (e.g., workload-independent noise)
  - **0**

# Inference of Complexity Transitions



Caller (**Order 0**, Constant)

RefreshListCtrl
**Order 0**

GetItemRelPath
**Order 1**

_listView.InsertItem
**Order 1**

...

...

...

Continue to search other children

Callee (**Order 1**, linear)
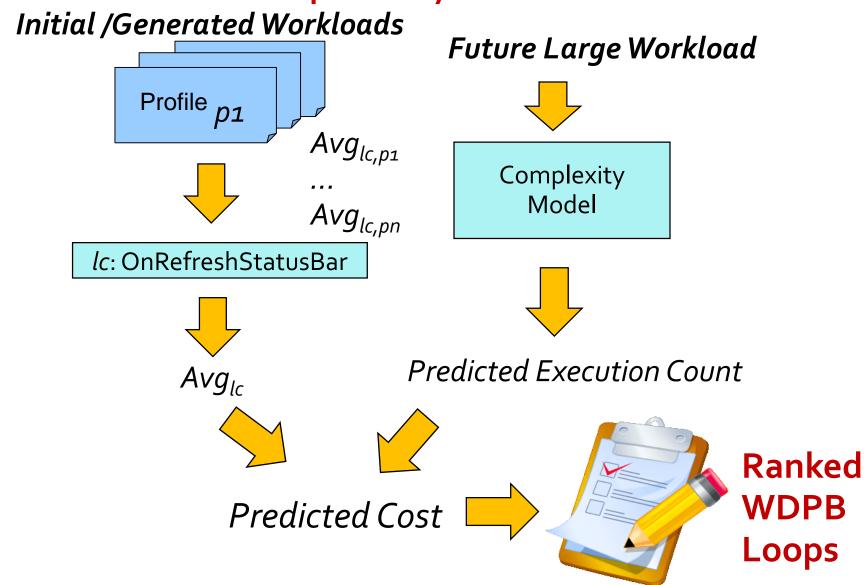
*Abstracted Model Comparison*

***Complexity Transition***

( RefreshListCtrl,    {GetItemRelPath, _listView.InsertItem, ...} )
**constant ➔ linear**

# Cost Prediction of Complexity Transitions

**Initial /Generated Workloads**

Profile $p_1$

$Avg_{lc,p_1}$
...
$Avg_{lc,pn}$

$lc$: OnRefreshStatusBar

$Avg_{lc}$

*Predicted Cost*

**Future Large Workload**

Complexity Model

*Predicted Execution Count*

**Ranked WDPB Loops**

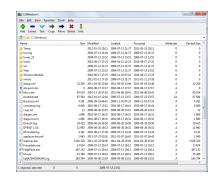# Evaluations of *DeltaInfer*

- Subjects : open-source GUI applications from SourceForge
  - 7-Zip: file manager       **7,280** LOC
  - Notepad++: text editor    **155,300** LOC



- RQs
  - RQ1: Effectiveness of WDPB Identification
  - RQ2: Effectiveness of Model Inference/Refinement
  - RQ3: Effectiveness of Context-Sensitive Analysis

# Evaluation Setup - Scenarios

|          | ID    | Scenario                                          | W. Param |
|----------|-------|---------------------------------------------------|----------|
| 7-Zip    | (S1)  | Open a folder                                     | # files  |
|          | (S2)  | Rename a file                                     | # files  |
|          | (S3)  | Select all items and then select the first item   | # files  |
|          | (S4)  | Create a folder                                   | # files  |
|          | (S5)  | Delete a file                                     | # files  |
| Notepad++| (S6)  | Open a file                                       | # lines  |
|          | (S7)  | Enter a character and save the file               | # lines  |
|          | (S8)  | Go to the last line                               | # lines  |
|          | (S9)  | Find a word not present in the file               | # chars  |
|          | (S10) | Cut and past the first character                  | # lines  |

# Evaluation Setup – Cont.

- Workload Selection
  - Representative Value Range (RVR)
    - Representative usage **[1, 1280]**
  - Validation Value Range (VVR)
    - Two times as RVR  **[1, 2560]**
  - Initial workloads
    - Initial workload groups: **{20,40,80}** , **{100,200,400}**
- Thresholds
  - Max refinement iterations: **20**
  - Threshold for $R^2$: **0.9**; Improvement threshold: **2%**
  - Prediction-error threshold: **5%**

# RQ1: WDPB Identification

- Manually inspect top-rank **complexity transitions**

- Report identified **performance bugs** for confirmation from developers

- Measure **cost coverage** of WDPBs as the workloads increase

# Example Bugs of 7-Zip

**Complexity Transition (constant to linear)**
(RefreshListCtrl, {GetItemRelPath, _listView.InsertItem, …}).

```
HRESULT RefreshListCtrl(…) {
…
    for (UInt32 i = 0; i < numItems; i++) {          WDPB loop
      …
        const UString relPath = GetItemRelPath(i);    Intensive temp-obj
      …                                               creation/destruction
      if (_listView.InsertItem(&item) == -1)
        return E_FAIL;
    }
    …
    listView.SortItems(CompareItems, (LPARAM) this);   Implicit
}                                                      WDPB
                                                       loop!
```
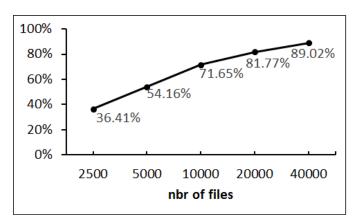
**Complexity Transition (constant to power-law)**
**(listviewProcedure, {CompareItems, …}).**

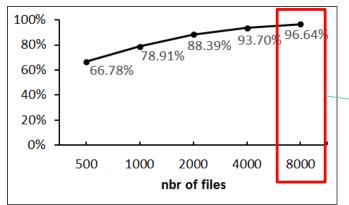**This bug is triggered in S1, S2, S3, S5**

# Example Bugs of Notepad++

**WDPB loop**

```
1  bool WrapLines(…)  {
2      …
3      while (lineToWrap < lastLineToWrap) {
4          if (WrapOneLine(surface, lineToWrap)) {
5              wrapOccurred = true;
6          }
7          lineToWrap++;
8      }
9      …
10  }
```

**Expensive computation**

**Complexity Transition (constant to linear)**
**(WrapLines, {WrapOneLine, …}).**

**This bug is triggered in S6, S7 S10**
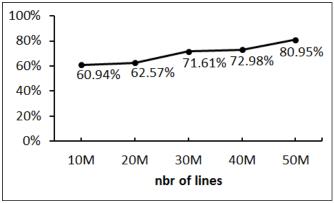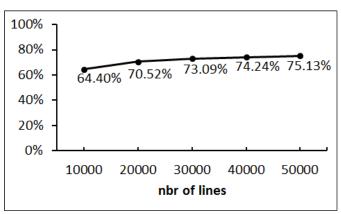
# Cost Coverage of WDPBs



7-Zip FM (S1: open a folder)



7-Zip FM (S3: select all and click any item)

**Nested loop!**



Notepad++ (S6: open a file)



Notepad++ (S9: Find a word)

- Identified WDPBs account for **75%+** of costs ➔ low probability of missing impactful WDPBs
- Cost coverage increases very fast for **nested loops**

# Bug Confirmations

- 7-Zip **(5 bugs)**
  - Bugs of **RefreshListCtrl** (S1, S2, S3, S5) confirmed with fix planned for next version
  - Bug of **RefreshStatusBar** (S4) introduced in release on Aug 05 and still remaining in latest release on Aug 11.
- Notepad++: **(5 bugs)**
  - Bug of **wraplines** (S6) found at the forum
  - Bugs caused by **wraplines** (S7,S8, S10), pending
  - New bug on **search for a not-present word** (S9), pending

# RQ2: Model Inference and Refinement

Initial Workloads

After Refinement

| ID | # Ite. | # WorkL | E. I. | E. E. |
|-----|--------|---------|--------|-------|
| (S1) | 4 | 4 | 35.95 | 0.62 |
| (S2) | 4 | 4 | 62.31 | 0.47 |
| (S3) | 4 | 4 | 29.68 | 0.85 |
| (S4) | 4 | 4 | 4.71 | 0.20 |
| (S5) | 3 | 3 | 5.94 | 0.18 |
| (S6) | 6 | 6 | 536.74 | 8.62 |
| (S7) | 5 | 5 | 455.00 | 7.69 |
| (S8) | 7 | 7 | 17.12 | 5.51 |
| (S9) | 4 | 4 | 138.36 | 1.83 |
| (S10) | 7 | 7 | 7.38 | 1.86 |

Avg relative errors of inferred complexity models

- **5 iterations (7 workloads)** to reach avg relative err 2.8%
- **Insensitive** to potential variations of initial workloads

# RQ2: Prediction Error of Cost

| ID | 10 (%) | 20 (%) | 50 (%) |
|---|---|---|---|
| (S1) | 3.18 | 4.45 | 6.16 |
| (S2) | 2.98 | 4.07 | 5.55 |
| (S3) | *1.40 | *1.60 | *1.86 |
| (S4) | 1.65 | 2.29 | 3.08 |
| (S5) | 1.58 | 2.19 | 2.95 |
| (Ave(7-Zip)) | *2.35 | *3.25 | *4.44 |
| (S6) | 18.51 | 6 | 47.24 |
| (S7) | 16.84 | 5 | 36.28 |
| (S8) | 16.80 | 7 | 35.23 |
| (S9) | 11.15 | 4 | 39.09 |
| (S10) | 10.79 | 7 | 24.63 |
| (Ave(Notepad++)) | 14.82 | 20.97 | 36.49 |

X RVR Upper Bound

**Developers optimize message processing during idle time**

Prediction error

- **4.4%** (7-Zip file manager): excluding S3
- **36.5%** (Notepad++ ): robust even under complex situations

# RQ3: Context Sensitivity

| ID | *DeltaInfer* | # L. | InSen. | # Missed by InSen |
|------|------|------|------|------|
| (S1) | 11 | 10 | 521 | 6 |
| (S2) | 21 | 19 | 579 | 12 |
| (S3) | 17 | 16 | 486 | 10 |
| (S4) | 21 | 19 | 640 | 12 |
| (S5) | 22 | 20 | 546 | 12 |
| (S6) | 10 | 10 | 509 | 3 |
| (S7) | 29 | 14 | 877 | 6 |
| (S8) | 10 | 10 | 526 | 5 |
| (S9) | 20 | 20 | 131 | 0 |
| (S10) | 12 | 11 | 861 | 3 |

**Real WDPB loops**

- Context helps **reduce false positives & negatives**
  - No context: **> 90%** of identified WDPB loops being **false positives**
  - No context: **40%** of DeltaInfer-identified WDPB loops being **missed**
- Context helps achieve only 14% of identified WDPB loops being **false positives** (top/low-level sys lib calls)

# Conclusion

- Predictive approach for WDPBs: *context-sensitive delta inference*
  - **Temporal** inference ➜ complexity models
    - **Deltas** of different **executions** (workloads)
  - **Spatial** inference ➜ complexity transitions
    - **Order deltas** of different **locations**

- Evaluations: effectively identifies **impactful WDPBs** (for causing **10** performance bugs)

# Thank You !

# Conclusion

- Predictive approach for WDPBs:
  ***context-sensitive delta inference***
  - **Temporal** inference ➜ complexity models
    - **Deltas** of different **executions** (workloads)
  - **Spatial** inference ➜ complexity transitions
    - **Order deltas** of different **locations**

- Evaluations: effectively identifies **impactful WDPBs** (for causing **10** performance bugs)

# Discussion

- Generalization to Other Types of Applications

- Multiple Workload Parameters

- Value-Dependent Performance Bottlenecks

- Scalability of Scenario-Based Profiling

# Software Performance

- An important quality of software
  - A kind of non-functional requirement
  - Characterized by the amount of work accomplished given time and resources

- Software performance matters
  - Software functionality and size grows faster than

  - been done to avoid related mista

Tablet

Data Center

# Two Significant Challenges in GUI Applications

- Complex contexts
  - GUI applications are event-driven applications
  - A program location may exhibit different complexities under different contexts

- Implicit WDPB loops
  - Event handlers can be invoked repetitively
    - E.g., selection change events for all items
  - No explicit loop statements
    - Cause challenges to manual inspection and static analysis

# Limitations of Traditional Approaches

- Traditional approaches
  - Performance testing (blackbox-random testing or manual testing)
  - Profiling (call-tree profiling and callstack sampling)

  > *41 out of 109 studied performance bugs are due to wrong assumption of workloads. [Jin et al. PLDI 2012]*

- Two major issues
  - Insufficiency
    - WDPBs may not surface on given workloads
    - Workload specifications are usually missing or outdated
  - Incompleteness:
    - WDPBs may overshadow other WDPBs

# Least-Squares Regression

- Linear Regression
  - Infers:  $y = A + Bw$
  - Minimizes: $Q(A,B) = \sum_{i=1}^{k}(y_{l_c,i} - (A + Bw_i))^2$
- Power-law Regression
  - Infers:  $y = Aw^B$
  - Minimizes: $Q(A,B) = \sum_{i=1}^{k}(y_{l_c,i} - (Aw_i^B))^2$
- How good does the model fit the data points?
  - Correlation coefficient
    - $R^2 = \dfrac{(\sum_{i=1}^{k} wy - k\overline{wy})^2}{(\sum_{i=1}^{k} w^2 - k\overline{w}^2)(\sum_{i=1}^{k} y^2 - k\overline{y}^2)}$,
    - $\overline{w}$ is the mean of w, $\overline{y}$ is the mean of y

# A Few Definitions

- Application *A*, location *l* and cost *y*
- Call graph G(E,V), calling context *c*, and execution profile *P*
- *k-profile Graph:* an annotated call graph, G(E, V ), where a location l with its corresponding vertex is annotated with a vector of counters for l on k workloads for each of its calling context c.

# Complexity Transitions

- A pair (n,M), such that:
  - 1. $n$ is a vertex (method) in the k-profile graph and $M$ is a subset of children vertices (callees) of $n$;
  - 2. $f_{n,c}(W)$ is the complexity model of n under the calling context $c$, and $f_{li,ci}(W)$ is the complexity model of the location $l_i$, where $l_i$ is a location in $M$ and the calling context $c_i$ is $c$ concatenated with $n$.
  - 3. $O(f_{li,}c_i(W))$ is at least 1 more than $O(f_{n,c}(W))$;
  - 4. $\forall l_i, l_j \in M$, i ≠ j, $O(f_{li,ci}(W)) = O(f_{lj,cj}(W))$.

# Model Inference and Refinement

```
 1: pc = −1 // previous model count
 2: e_pre = −1 // previous error
 3: for ite = 0; ite < max; ite = ite + 1 do
 4:     kG = AlignProfiles(P)
 5:     x = GetWorkloads(P)
 6:     M = RegressionLearning(x, kG)
 7:     e_M = {}
 8:     for all vp in VP do
 9:         w = GetWorkload(vp)
10:         e_vp = 0.0
11:         for all m in M do
12:             r_w = Predict(w, m)
13:             a_m = GetActual(vp, m)
14:             e_vp = e_vp + Abs(a_m − r_w) / a_m
15:         end for
16:         e_M = e_M.Add(e_vp / M.Count)
17:     end for
18:     e_total = Sum(e_M) / VP.Count
19:     if e_pre − e_total < threshold_imp then
20:         break // improvement is below the threshold
21:     end if
22:     if e_total < threshold_e AND pc == M.Count then
23:         break // accuracy is acceptable
24:     else
25:         np = NewWorkload(P, VP, e_M)
26:         P = P.Add(np)
27:         pc = M.Count
28:         e_pre = e_total
29:     end if
30: end for
31: return  M
```

Align Profiles
- Align locations using calling contexts
- Extract execution vector for each location under each calling context

Regression Learning

Model Validation

Termination Checks

Select New Workloads
- Assumption: a new workload at the area with the highest prediction error improves most

# Cost Prediction of Complexity Transitions

- Compute $avg_{lc,p}$ for each location $l_c$ on each profile $p$
  - E.g., for $p_1$, $Cost(refreshList_c) = 1s$, $ExeCount(refreshList_c) = 100$, $avg_{lc,p} = 1/100$ s

- Compute $avg_{lc} = average(avg_{lc,p})$

- Given a workload value $w$, $pred_{lc,w} = f_{lc}(w)$

- Get $Cost_{lc,w} = pred_{lc,w} * avg_{lc}$