VIRGINIA TECH™

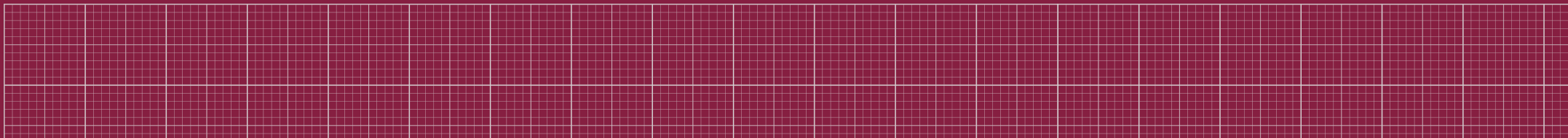**D-Goldilocks:**
**Automatic Redistribution of**
**Remote Functionalities**
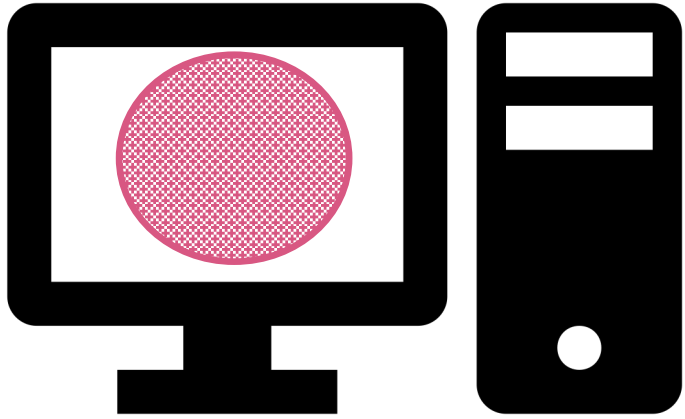**for Performance and Efficiency**
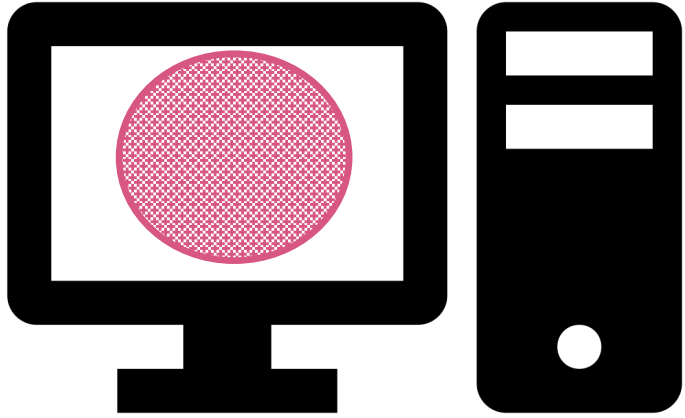
KIJIN AN AND ELI TILEVICH
SOFTWARE INNOVATIONS LAB

# Distribution

- **Why Distributed Computing?**
  - Take advantage of remote **computing resources**
  - Improve **performance** and/or **efficiency**
- **Distribution Benefits**
  - Access superior remote resources
  - Share the **computational load**
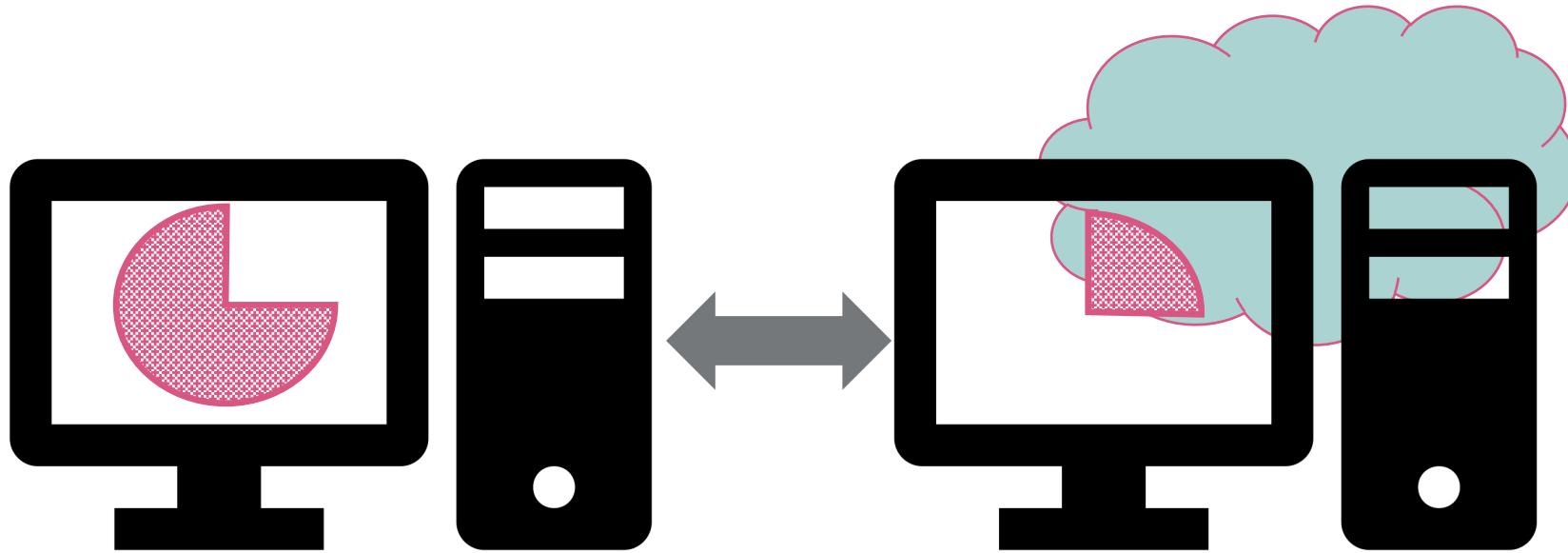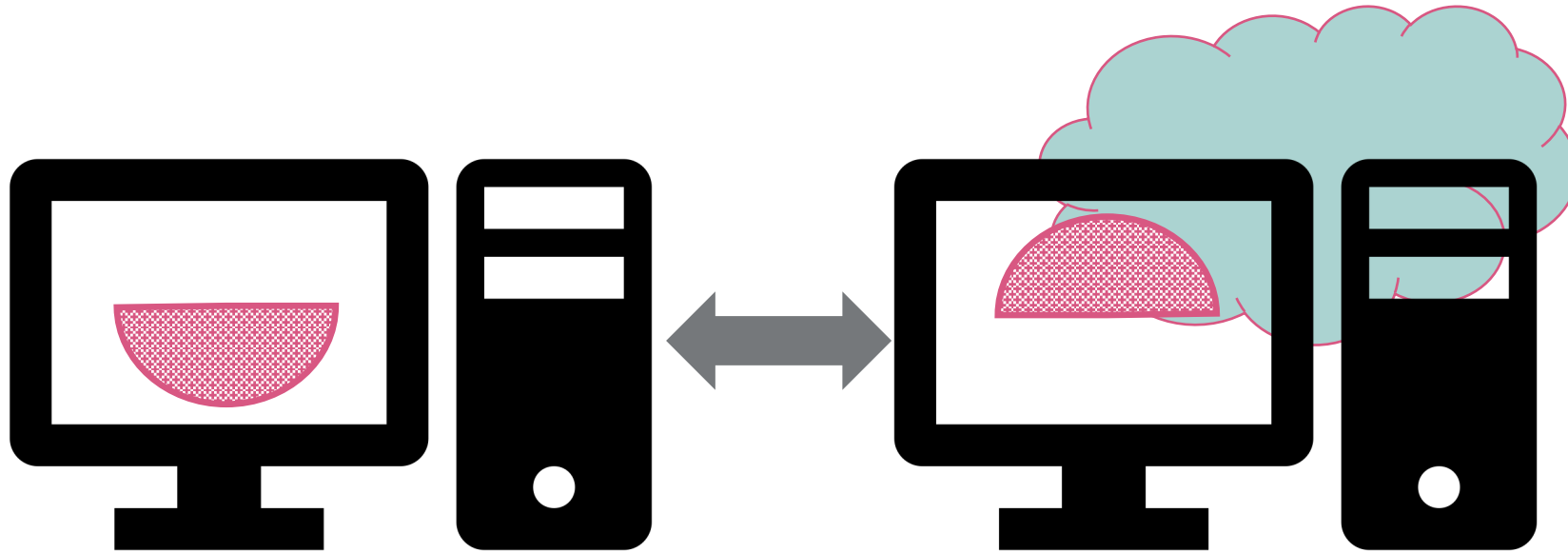- **Distribution Costs**
  - Communication Overhead
  - Partial failure
  - Security

# Distribution

# Distribution

# Distribution

# Distribution

# Distribution

# Distribution (Granularity of Remote Service)

**Too Crude!**
(Granularity of Remote Service)

# Distribution (Granularity of Remote Service)

(Asynchronous Executions)

**Too Fine!**
(Granularity of Remote Service)

# Too Much Remote Execution not always beneficial
## :Nano Service Anti-pattern[Moha 2012 et. al]

*Slower Speed*

*Low Overhead*

*Bulk*

**(Granularity of Remote Service)**

*Faster Speed*

*High Overhead*

*Nano*

*"Performance"*        *"Efficiency"*

# Motivating Real-world's Example: Bookworm

**Client-Side**

**Server-Side**

`/api/ladydog`

Executed Sequentially

GetSentence→
getUniqueVocaburary→
GetDialog→getColon→
getSemiColon

**Initial** distribution

# Motivating Real-world's Example: Bookworm

**Client-Side**

**Server-Side**
`/api/ladydog`

Executed Sequentially

GetSentence→
getUniqueVocaburary→
GetDialog→getColon→
getSemiColon

"No Results Until All tasks to complete!"

**Execute in Bulk!**

**Initial** distribution

# Motivating Real-world's Example: Bookworm

**Client-Side**



**Server-Side**

`/api/ladydog`

GetSentence

getUniqueVocaburary

GetDialog

getColon

getSemiColon

**Independent of each other**

# Motivating Real-world's Example: Bookworm

**Client-Side**

**Server-Side**

`/api/ladydog`

GetSentence

getUnique Vocaburary

GetDialog

getColon

getSemiColon

**Splitting into smallest Units to invoke Remotely**

# Motivating Real-world's Example: Bookworm

**Client-Side**

**Server-Side**
`/api/ladydog`

**Too much Distribution**

GetSentence

getUniqueVocaburary

GetDialog

getColon

getSemiColon

**"It's Faster Execution!"**
By invoking all together, "Asynchronously"

# Goldilocks Principle

Too Crude                                    Too Fine

# Goldilocks Principle

Just  the right one!

Too Crude                                                                    Too Fine

# D-Goldilocks

**Client-Side**

**Server-Side** (`/api/ladydog`)

**Too Small Distribution**

getVoc    GetC

GetSen    GetDial

getCol

"Too Crude"
(Long **Latency**)

*Bookworm* Data insights into classic stories

The Dead
James Joyce

The Cast of
Amontillado
Edgar Allan Poe

with
Dog

The
Red
Roo

## The Lady with the Pet Dog

by Anton Chekhov

Level of granularity

**Too Much Distribution**

Commutation
Overheads

getVoc    GetC

GetSen    GetDial

getCol

"Too Fine"
(Too much **Overheads**)

# D-Goldilocks

**Client-Side**

**Server-Side** (`/api/ladydog`)

**Too Small Distribution**

getVoc   GetC
GetSen   GetDial
getCol

"Too Crude"
(Long **Latency**)

Bookworm — Data insights into classic stories

The Dead — James Joyce
The Cast of Amontillado — Edgar Allan Poe

**"Right"** Redistribution

(Right boundary)

getVoc

GetC
GetDial

GetSen
getCol

Level of granularity

The Lady with the Pet Dog
by Anton Chekhov

Goldilocks Principle

**Too Much** Distribution

Commutation Overheads

getVoc   GetC

GetSen   GetDial

getCol

"Too Fine"
(Too much **Overheads**)

# Problem Formulation

- Determine which functional distribution from the client's standpoint would minimize the <mark>cost of distributed execution</mark>

$$C_{Dist\_Exec}(\mathbf{r}) = \alpha \cdot latency(\mathbf{r}) + (1-\alpha) \cdot \Sigma \ resource(\mathbf{r})$$

Execution Time
(Performance)

Consumed Resource
(Efficiency)

Normalizing
Parameter

# Problem Solution Outline

- Redistribution operations:
  - Partition
    - $[r_1,...,r_k] = \textbf{partition}(r)$

  - Batch
    - $r_h = \textbf{batch}([r_1,...,r_n])$

(A Remote Execution)  (Independently Invocable)

$r$  →  *partition*  →  $r_1$  $r_2$  $r_{k-1}$  $r_k$

(Multiple
Remote Executions)  (Greater granularity)

$r_1$  $r_{n-1}$  $r_n$  →  *batch*  →  $r_h$

# How to restructure Remote Services?

- **Client Insourcing Refactoring** [WWW '20]
  - Undoing Distribution

# How to restructure Remote Services?

- **Client Insourcing Refactoring** [WWW '20]
  - Undoing Distribution

# How to restructure Remote Services?

- **Client Insourcing Refactoring** [WWW '20]
  - Undoing Distribution

# How to restructure Remote Services?

- **Client Insourcing Refactoring** [WWW '20]



① Identifying Entry/Exit Points of Remote Functionality

② Constraints Solving & Code Transformation

# How to restructure Remote Services?

- **Client Insourcing Refactoring as re-distribution framework**

(Original Distribution)          (Re-Distribution)

$r$

$r'$

Remote

Local

**Client Insourcing Refactoring** [*WWW '20*]

r_local ----> r_local'

[Any Refactoring]
For Centralized Apps

[Any Distributing
Frameworks]

[Kwon ICDCS'13,
EXTREMEJS '12,..]

# Restructuring: Partition

(Original Distribution)

(Re-Distribution)

r

*partition*

$r_1$

$r_2$

$r_{k-1}$

$r_k$

**Remote**

*(Equivalent Variant)*

*(Independent sets of functions)*

**Local**

*Client Insourcing Refactoring*

r_local

$r\_local_1$

$r\_local_2$

...

$r\_local_{k-1}$

$r\_local_k$

**Partitioning**

# Restructuring: Partition

(Original Distribution)

(Re-Distribution)

**r**

$partition$

$r_1$

$r_2$

Maximum # of Distributions

$r_{k-1}$

$r_k$

Remote

*(Equivalent Variant)*

*(Independent sets of functions)*

Local

*Client Insourcing Refactoring*

r_local

$r\_local_1$

$r\_local_2$

...

$r\_local_{k-1}$

$r\_local_k$

**Partitioning**

**Distribution**

# Restructuring: Partition

(Original Distribution)

(Re-Distribution)

**r**

$partition$

$r_1$

$r_2$

$r_{k-1}$

$r_k$

Remote

Local

(Equivalent Variant)

(Independent sets of functions)

*Client Insourcing Refactoring*

r_local

ALL function decls

$r\_local_1$

$r\_local_2$

...

$r\_local_{k-1}$

$r\_local_k$

**Partitioning**

Partioning **r_local** into independent sets of functions
- Initial Candidates: **ALL function decls** in **r_local**
- Find partitions that are independent each other by using **Dependency analysis** for "Control flows" and "global variables" between **function decls**

# Restructuring: Partition & Batch



(Original Distribution)

(Re-Distribution)

$r$ → *partition* → $r_1$ $r_2$ $r_{k-1}$ $r_k$ → *batch\** → $r_1$ $r_{k-1}$ $r_k$ $r_2$ ...

*(Equivalent Variant)*

Remote

Local

*Client Insourcing*

r_local → r_local$_1$ r_local$_2$ ... r_local$_{k-1}$ r_local$_k$ → r_local$_1$ ... r_local$_{k-1}$ r_local$_k$ r_local$_2$

*(Inline Function Refactoring)*

**Partitioning**          **Batching**

30

# Restructuring: Partition & Batch



(Original Distribution)

(Re-Distribution)

$r$

$partition$

$r_1$ $r_2$ $r_{k-1}$ $r_k$

$batch*$

$r_1$ $r_{k-1}$ ... $r_k$ $r_2$

Remote

Local

(Equivalent Variant)

Client Insourcing

r_local

Partitioning

$r\_local_1$ $r\_local_2$ ... $r\_local_{k-1}$ $r\_local_k$

(Inline Function Refactoring)

$r\_local_{1 ...}$ $r\_local_{k-1}$ $r\_local_k$ $r\_local_2$

How to distribute? (local->remote)

# Restructuring: Partition & Batch



(Original Distribution)

(Re-Distribution)

$r$

$partition$

$r_1$   $r_2$   $r_{k-1}$   $r_k$

$batch*$

...

$r_1$   $r_{k-1}$   $r_k$   $r_2$

Remote

Local

*(Equivalent Variant)*

*Client Insourcing*

r_local

r_local$_1$
r_local$_2$
...
r_local$_{k-1}$
r_local$_k$

r_local$_{1\,...}$
r_local$_{k-1}$
r_local$_k$

r_local$_2$

Remote Façade[1]

Remote Façade[h]

**(Inline Function Refactoring)**

**Partitioning**

**Batching Remote Invocation**

# Batching Remote Invocation(Batch)

- **Distributing Programing Pattern** [Fowler '02, Ibrahim et.al ECOOP '09]

f1(p1)

f2(p1)

f3(p2)

**Client DTO**
f1_f2_f3

f1_f2_f3(p1, p2, p3)

DTO_script

**Remote Façade**
f1_f2_f3

f1_f2_f3.f1(p1)

f1_f2_f3.f2(p1)

f1_f2_f3.f3(p2)

The **Client DTO** stub accumulates the fine-grained service invocations then, transfers them (parameters) in bulk.

The **remote Façade** function sequentially invokes the bundled services. Then, it combines their execution results and returns in bulk.

# Process for D-GOLDILOCKS



Initial Distribution
(Full-Stack JS app)

Centralized
Variant function

Façade

BRI patterns

DTO

Distribution
Middleware
(multi-cored)

(Headless
Browser Testing framework)

Client
Insourcing
(z3 Solver)

Partitioning

Independent
Sub-function

Batching
Remote
Invocation

templates

new
Distributions

[[0], [1], [6, 8], [2, 3], [5, 7]],
[[0], [1], [2], [3], [5, 7], [6, 8]]
[[2], [3], [6], [8], [0, 1], [5, 7]],
[[0], [3], [6], [8], [1, 2], [5, 7]]
[[0], [1], [2], [6], [8], [3, 5, 7]]
[[0], [1], [2], [3], [6], [5, 7, 8]]
[[0], [1], [6], [7], [2, 3], [5, 8]],
[[0,1], [2,3], [5, 8], [6, 7]] ,

Batching Parameter Sets

Measuring
Cost Functions

# Process for D-GOLDILOCKS



Initial Distribution
(Full-Stack JS app)

JS

JS

Centralized
Variant function

**Client
Insourcing**
(z3 Solver)

**Partitioning**

Independent
Sub-function

Façade

BRI patterns

DTO

**Distribution
Middleware**
(multi-cored)

templates

**Batching**
Remote
Invocation

(Headless
Browser Testing framework)

new
Distributions

Measuring
Cost Functions

[[0], [1], [6, 8], [2, 3], [5, 7]] ,
[[0], [1], [2], [3], [5, 7], [6, 8]]
[[2], [3], [6], [8], [0, 1], [5, 7]]
[[0], [3], [6], [8], [1, 2], [5, 7]]
[[0], [1], [2], [6], [8], [3, 5, 7]]
[[0], [1], [2], [3], [6], [5, 7, 8]]
[[0], [1], [6], [7], [2, 3], [5, 8]]
[[0,1], [2,3], [5, 8], [6, 7]] ,

Batching Parameter Sets

# Process for D-GOLDILOCKS



Initial Distribution
(Full-Stack JS app)

Centralized Variant function

Client Insourcing (z3 Solver)

Partitioning

Independent Sub-function

Façade

BRI patterns

DTO

Distribution Middleware (multi-cored)

templates

Batching Remote Invocation

(Headless Browser Testing framework)

new Distributions

Measuring Cost Functions

[[0], [1], [6, 8], [2, 3], [5, 7]] ,
[[0], [1], [2], [3], [5, 7], [6, 8]]
[[2], [3], [6], [8], [0, 1], [5, 7]]
[[0], [3], [6], [8], [1, 2], [5, 7]]
[[0], [1], [2], [6], [8], [3, 5, 7]]
[[0], [1], [2], [3], [6], [5, 7, 8]]
[[0], [1], [6], [7], [2, 3], [5, 8]]
[[0,1] , [2,3] , [5, 8] , [6, 7]] ,

Batching Parameter Sets

# Process for D-GOLDILOCKS

Façade

BRI patterns

Distribution
Middleware
(multi-cored)

DTO

(Headless
Browser Testing framework)

Initial Distribution
(Full-Stack JS app)

Centralized
Variant function

JS

JS

templates

new
Distributions

Client
Insourcing
(z3 Solver)

Partitioning

Independent
Sub-function

Batching
Remote
Invocation

Measuring
Cost Functions

[[0], [1], [6, 8], [2, 3], [5, 7]] ,
[[0], [1], [2], [3], [5, 7], [6, 8]]
[[2], [3], [6], [8], [0, 1], [5, 7]]
[[0], [3], [6], [8], [1, 2], [5, 7]]
[[0], [1], [2], [6], [8], [3, 5, 7]]
[[0], [1], [2], [3], [6], [5, 7, 8]]
[[0], [1], [6], [7], [2, 3], [5, 8]]
[[0,1] , [2,3] , [5,8] , [6, 7] ,

Batching Parameter Sets

- **Instrumenting Time: a couple of MINs ~ a couple of HOURs** for each subject (DELL-OPTIPLEX5050)
- Depending on possible Combinations to batch them

# **Evaluation**: Research Questions

- **RQ1:—Value:** How much **programmer effort** is **saved** by D-GOLDILOCKS's automatic redistribution operations?
- **RQ2:—Cost Model Correctness:** How applying the partition and batch operations affect the distributed execution's "**latency**" and "**consumed resources**" attributes?
- **RQ3:—Utility of Cost Model for Redistribution:** How useful is the **cost function** for guiding redistribution decisions?
- **RQ4:—Energy Consumption:** What is the effect of redistribution on the amount of **energy consumed** by the client?

# Subject Full-stack JavaScript Apps :Original Performance and Efficiency

Original "Latency"     Original "CPU utilization"(Resource)

| Remote Services | L(ms) | ΣTCPU | $f_{CI}^{LOC}$ | $f_{decl}$ | $f_{sub}^{ind}$ | IDI |
|---|---|---|---|---|---|---|
| /api/ladypet | 77.83 | 337 | 394 | 9 | 8 | 1.6M |
| /api/thedea | 164.62 | 695 | 394 | 9 | 8 | 1.6M |
| /api/bigtrip | 42.11 | 304 | 394 | 9 | 8 | 1.6M |
| … (Total 12 Subjects from 4 Full-Stack Apps) | | | | | | |
| /string-fasta | 29.85 | 328 | 38 | 5 | 2 | 76 |
| /cflow-rec | 35.43 | 326 | 49 | 4 | 3 | 245 |
| /prprty/brokers | 20.64 | 323 | 379 | 3 | 3 | 1.5K |

# Subject Full-stack JavaScript Apps

:How many Lines of Code($f_{CI}^{LOC}$) and Independent sub-functions are in the original remote functionality(**Centralized Variant**)?

| Remote Services | L(ms) | ΣTCPU | $f_{CI}^{LOC}$ | $f_{decl}$ | $f_{sub}^{ind}$ | IDI |
|---|---|---|---|---|---|---|
| /api/ladypet | 77.83 | 337 | 394 | 9 | 8 | 1.6M |
| /api/thedea | 164.62 | 695 | 394 | 9 | 8 | 1.6M |
| /api/bigtrip | 42.11 | 304 | 394 | 9 | 8 | 1.6M |
| **... (Total 12 Subjects from 4 Full-Stack Apps)** | | | | | | |
| /string-fasta | 29.85 | 328 | 38 | 5 | 2 | 76 |
| /cflow-rec | 35.43 | 326 | 49 | 4 | 3 | 245 |
| /prprty/brokers | 20.64 | 323 | 379 | 3 | 3 | 1.5K |

# Subject Full-stack JavaScript Apps

Initial Distribution



(Headless Browser Testing framework)

Templates

**Client Insouricing& Partitioning**

Independent Sub-function

**Batching** Remote Invocation

$f_{CI}^{LOC}$

new Distributions

```
[[0], [1], [6, 8], [2, 3], [5, 7]] ,
[[0], [1], [2], [3], [5, 7], [6, 8]] ,
[[2], [3], [6], [8], [0, 1], [5, 7]] ,
[[0], [3], [6], [8], [1, 2], [5, 7]] ,
[[0], [1], [2], [6], [8], [3, 5, 7]] ,
[[0], [1], [2], [3], [6], [5, 7, 8]] ,
[[0]
[[0]
```

All combinations of Batching of $f_{sub}^{ind}$

41

| $f_{CI}^{LOC}$ | $f_{decl}$ | $f_{sub}^{ind}$ | IDI |
|---|---|---|---|
| 394 | 9 | 8 | 1.6M |
| 394 | 9 | 8 | 1.6M |
| 394 | 9 | 8 | 1.6M |
| .... | | | |
| 38 | 5 | 2 | 76 |
| 49 | 4 | 3 | 245 |
| 379 | 3 | 3 | 1.5K |

# Subject Full-stack JavaScript Apps

:RQ1 Value : How much **Programmer Effort** is **Saved** by D-GOLDILOCKS's automatic redistribution operations?

Initial Distribution

Templates

(Headless Browser Testing framework)

Client Insouricing& Partitioning

Independent Sub-function

$f_{CI}^{LOC}$

**Batching** Remote Invocation

new Distril

IDI

[[0], [1], [6, 8], [2, 3], [5, 7]] ,
[[0], [1], [2], [3], [5, 7], [6, 8]] ,
[[2], [3], [6], [8], [0, 1], [5, 7]] ,
[[0], [3], [6], [8], [1, 2], [5, 7]] ,
[[0], [1], [2], [6], [8], [3, 5, 7]] ,
[[0], [1], [2], [3], [6], [5, 7, 8]] ,
[[0
[[0

All combinations of Batching of $f_{sub}^{ind}$

394 × 4139
~= 1.6 × 10⁶
ULOCs

| $f_{CI}^{LOC}$ | $f_{decl}$ | $f_{sub}^{ind}$ | IDI |
|---|---|---|---|
| 394 | 9 | 8 | **1.6M** |
| 394 | 9 | 8 | **1.6M** |
| 394 | 9 | 8 | **1.6M** |
| .... | | | |
| 38 | 5 | 2 | **76** |
| 49 | 4 | 3 | **245** |
| 379 | 3 | 3 | **1.5K** |

# Latency

Latency[ms]



# of Remote Executions

- **RQ2: Model Correctness(latency)**
- The larger the number of new remote functionalities, the smaller is the aggregate average latency $(\text{latency}(r) = \frac{1}{n}\sum_{i=1}^{n} T(r_i))$
- Splitting a single long-running remote function into a small number of asynchronously invoked parts



Asynchronously
Invoking parts

Latency

/api/thereadroom
/api/thebigtrip
/api/ladypet
/api/the_d
/api/thegift
/api/thecask
/api/wallpaper
/api/offshore

# Resource(CPU Utilization)

$\Sigma T_{CPU}$



# of Remote Executions

**RQ2: Model Correctness (Resources)**

- We measured total CPU Utilization to invoke a remote service r: resource(r) $= \sum_{i=1}^{n} CPU(r_i)$

- **Consuming Client's Resource a lot to invoke multiple remote executions, propositionally to # of invocations**

Latency(ms)

# of Remote Invocations

Sigma($T_{CPU}$)

# of Remote Invocations

lizing factor for the
rms by scaling the observed
nts $\alpha = \overline{L}/\overline{\Sigma T_{cpu}} = \mathbf{0.9281}$

$\mathbf{) \cdot \Sigma resource(r)}$

Consumed Resource
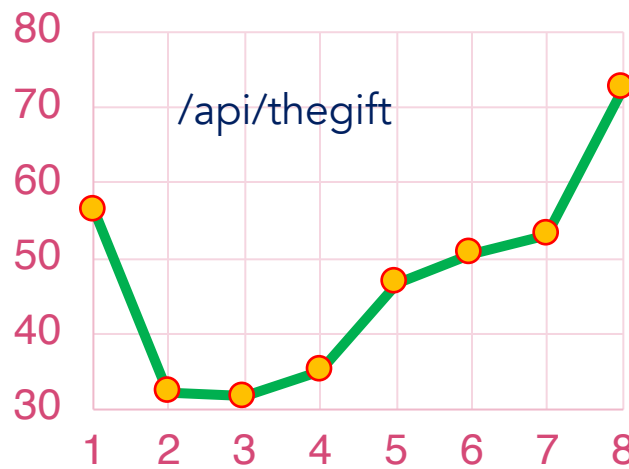(Efficiency)

Normalizing factor

Latency(ms)

# of Remote Invocations

Sigma($T_{CPU}$)

# of Remote Invocations

47

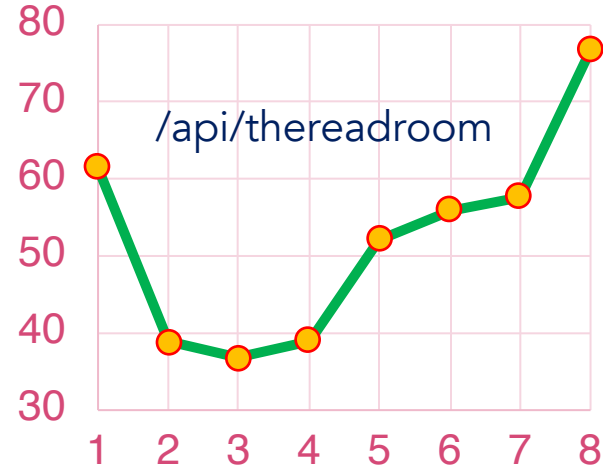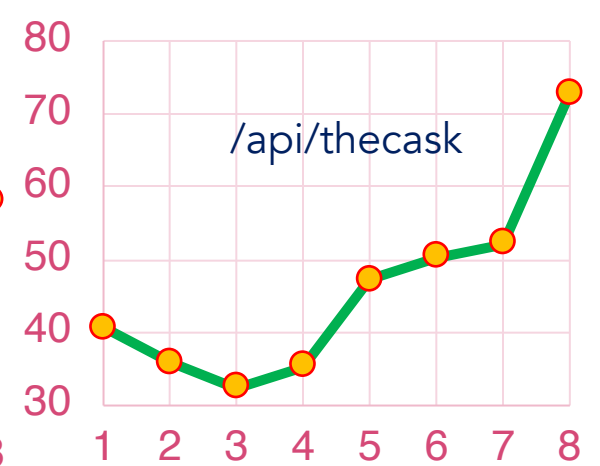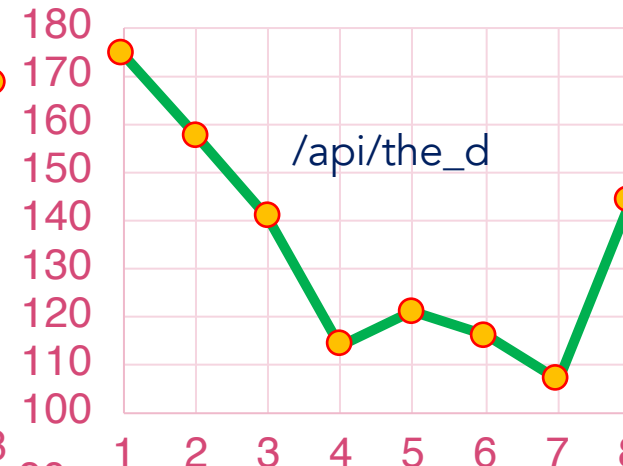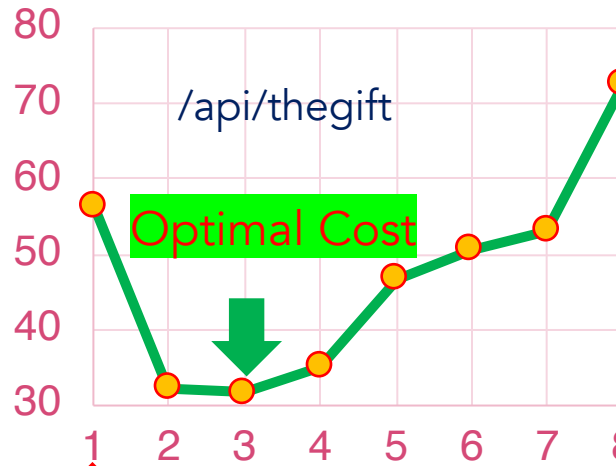# Cost Function

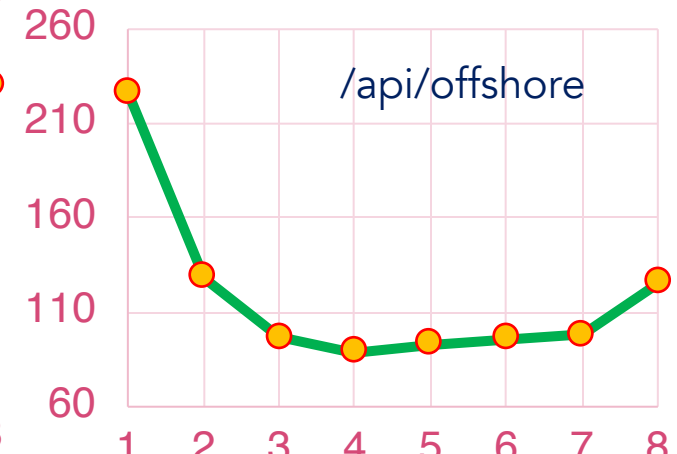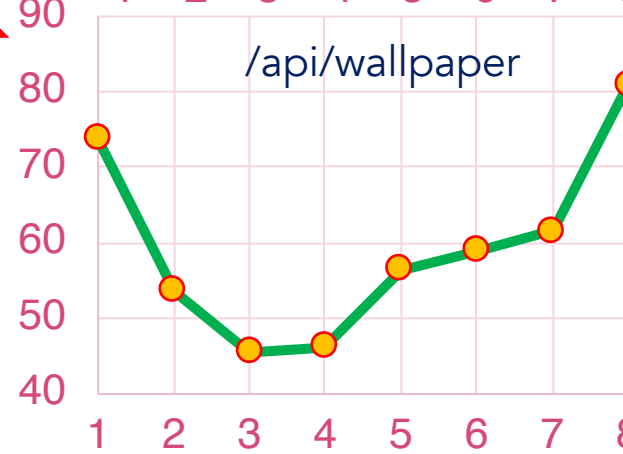**RQ3:—Utility of Cost Model**

# Cost Function

**RQ3:—Utility of Cost Model**



Too **Small** Distribution

Too **Much** Distribution

49

# Energy Consumption

- **RQ4:** What is the effect of redistribution on the amount of energy consumed by the client?

  - We natively build the subject app **(BookWorm)** by using **Apache Cordova**

  - **PowerTutor** [L Zhang et.al]: a model-based energy profiler for mobile apps

  - **Energy Consumptions** (EC) for <u>Original</u>, <u>Worse</u>, and <u>Best</u>

Subject App
*.HTML,
*.JavaScript,
*.CSS

**Apache Cordova**          **QISKIW-L24-HUAWEI**          **PowerTutor**
                            **(Marshmallow)**

50

# Energy Consumption(Original)

- Cost Function Versus. Energy Consumption

| | | |
|---|---|---|
| $EC_{original\_dist}$ | Original | 8.4mJ |
| $EC_{best\_dist}$ | MIN Cost | 13.4mJ |
| $EC_{worst\_dist}$ | MAX Cost | 47.4mJ |

Lowest Energy Consume But Highest Latency!

/api/thereadroom

getSentence

getDialog

getPeriods

getVocabulary

getQuestions

getColons

getCommas

getSemiColons

# Energy Consumption(Worst)

- Cost Function Versus. Energy Consumption

| | | |
|---|---|---|
| $EC_{original\_dist}$ | Original | 8.4mJ |
| $EC_{best\_dist}$ | MIN Cost | 13.4mJ |
| $EC_{worst\_dist}$ | MAX Cost | 47.4mJ |

**600%** More Energy Consumption!

/api/thereadroom

getSentence

getDialog

getPeriods

getVocabulary

getQuestions

getColons

getCommas

getSemiColons

# Energy Consumption(Best Dist. by D-Goldilocks)

- Cost Function Versus. Energy Consumption

| | | |
|---|---|---|
| $EC_{original\_dist}$ | Original | 8.4mJ |
| $EC_{best\_dist}$ | MIN Cost | 13.4mJ |
| $EC_{worst\_dist}$ | MAX Cost | 47.4mJ |

**50%** More Energy Consumption

getSentence→
→getDialog→
getPeriods

getVocabulary

getQuestions→
→getColons→
getCommas→
getSemiColons

53

# Conclusion

- A set of domain-specific **automatic refactorings** for **reshaping** and **redistribution**.

- A **cost function-based** heuristic for identifying how to improve the **performance** and **efficiency** of distributed apps by reshaping the original distribution, which was **too crude**.

- A systematic evaluation of our approach's **value**, **utility**, and **efficiency** for our reference implementation **"D-Golilocks"**
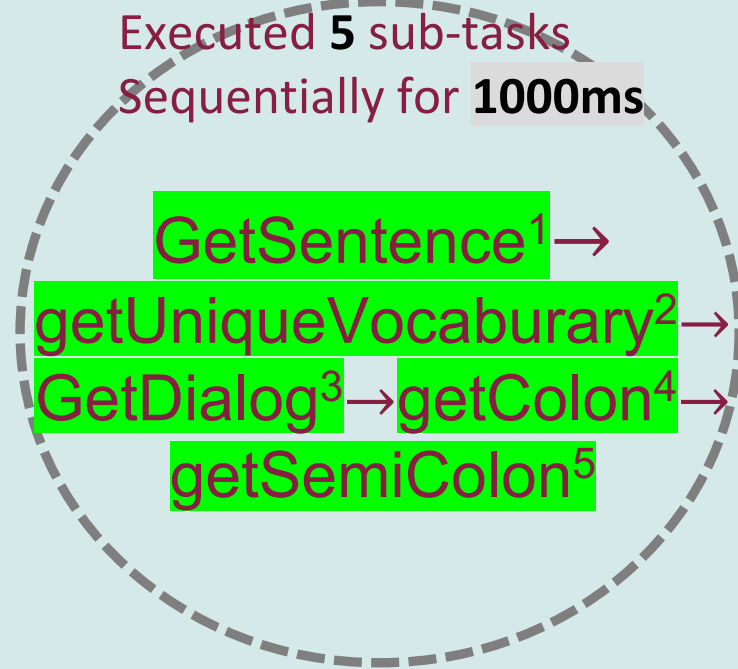
# Future Work

- Problem Formulation & Solution for different **Capacity** of **Clients** and **Servers** including **Network condition**

- Adaptation to **Edge Computing** for addressing their resource constraints and execution volatility

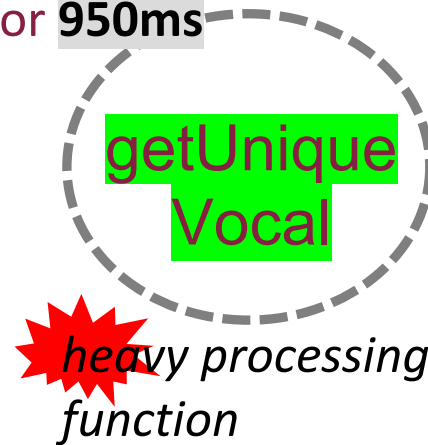- Other types of Software **Evolution Scenarios**

Q n A

Thank you!

# **Appendix-1**: Can multiple Executions reduce the aggregate average latency non-linearly?: $latency(r) = \frac{1}{n}\sum_{i=1}^{n} T(r_i)$
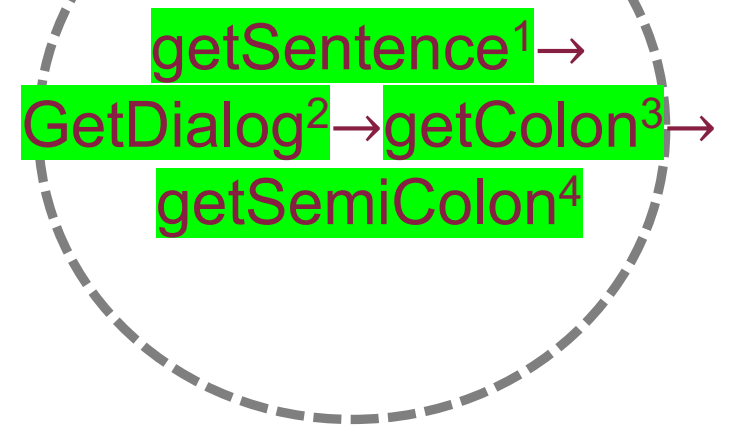
Executed **5** sub-tasks
Sequentially for **1000ms**

GetSentence[1]$\to$
getUniqueVocaburary[2]$\to$
GetDialog[3]$\to$getColon[4]$\to$
getSemiColon[5]

Executed a heavy
processing function
for **950ms**

getUnique
Vocal

*heavy processing
function*

Executed **4** sub-tasks
Sequentially for **150ms**

getSentence[1]$\to$
GetDialog[2]$\to$getColon[3]$\to$
getSemiColon[4]

**Latency_1**
= 1/5 (1000+1000+
1000+1000+1000) = **1000 (ms)**

**Latency_2**
= 1/5 (950+150+
150+150+150) = **310 (ms)**

# **Appendix-2**: Code-Base Example

```
//original Server:server.js
function getSenAvg(array){...};
function getVoca(str){...};
...
app.get('/api/ladypet',
    function(req, res){...});
```

```
//original Client:app.js
$scope.getLadyWithPetDog =
    function() {...
  $http.get('/api/ladypet').then(
      function(response){
    var text = response.data; ..
  });/*remote invocation*/ }
```

```
//after Client Insourcing:app.js
//Insourced remote functions
function getSenAvg(array){...};
function getVoca(str){...};
...
function ladypet_local(){
//invoke every subtasks
...};
$scope.getLadyWithPetDog =
    function() {...
    //from remote to local
    var text = ladypet_local();
...}
```

```
//after Redistribution:index.html
<!DOCTYPE html>
<script src="./app.js">
 ...
 ClientDTO.b_param = BATCH_PARAM;
//Batched Invocations:
 getSenAvg=ClientDTO(getSenAvg);
 getVoca=ClientDTO(getVoca);
 ...
</script>
```

Client
Insourcing

Batching fine-grained
service invocations

58