

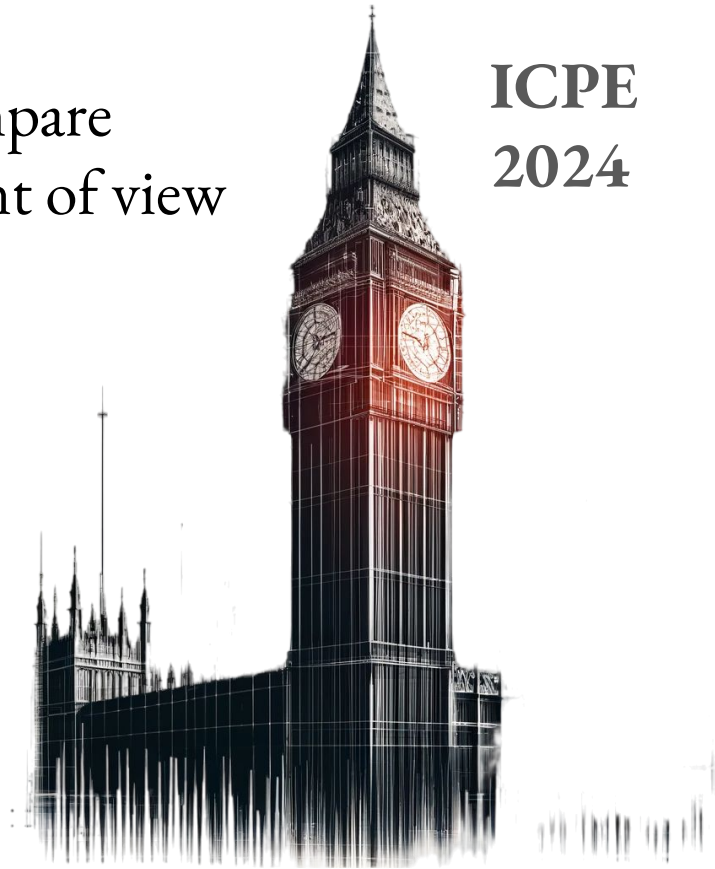
Distributed Trace Compare

Technical and practical point of view

ICPE
2024

Maryam Ekhlesi
Nasser Ezzati-Jivan

May 2024



Who we are?

Nasser Ezzati-Jivan

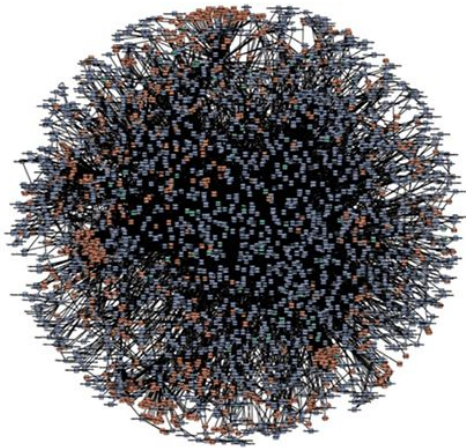
- From Brock University
- Assistant Professor
- Professor, software engineer, and team leader since 2008

Maryam Ekhlesi

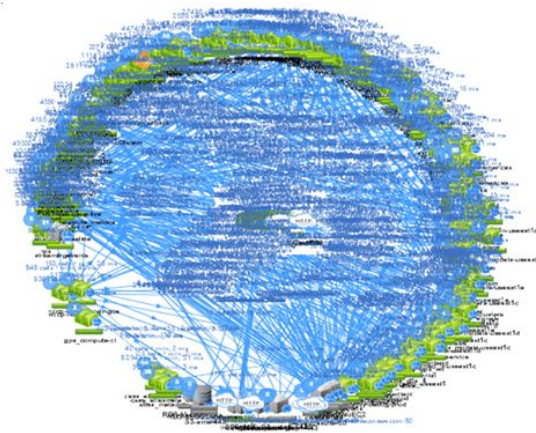
- From Polytechnique Montreal University, DORSAL Lab.
- Ph.D. Candidate
- Software Performance researcher
- Software engineer since 2011



Goal



amazon.com



NETFLIX

Agenda

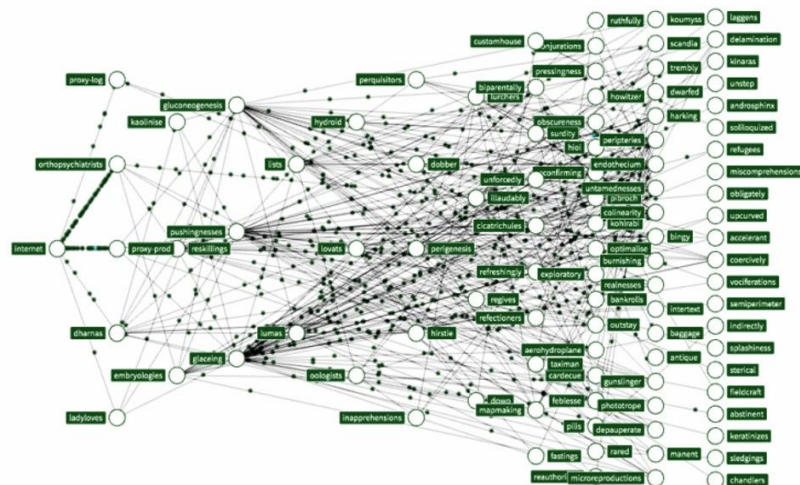
1. Introduction (5 min)
2. Motivations (5 min)
3. Basic concepts (60 min)
 - Low-level Tracing, End-to-end Tracing, Debugging, Profiling, Logging, Monitoring
 - Kernel/user space
 - Process/ thread states
1. Detecting/Diagnosing and fixing real performance problems (60 min)
 - Introducing real performance problems
 - Introducing Differential flame graph
 - Detecting performance problems
 - Diagnosing performance problems
2. Hands-on exercises (30 min)
 - Installing TraceCompass
 - Collecting Traces with LTTng
3. Conclusion (5 min)
4. Questions and Answers (5 min)

Special Thanks!

- Professor Michel Dagenais
 - Full Professor at Polytechnique Montreal University
- Fatemeh Faraji Daneshgar
 - Research Associate at Polytechnique Montreal University
 - Developer of Differential Flame Graph in TraceCompas
- TraceCompass Team
- EfficiOS Team

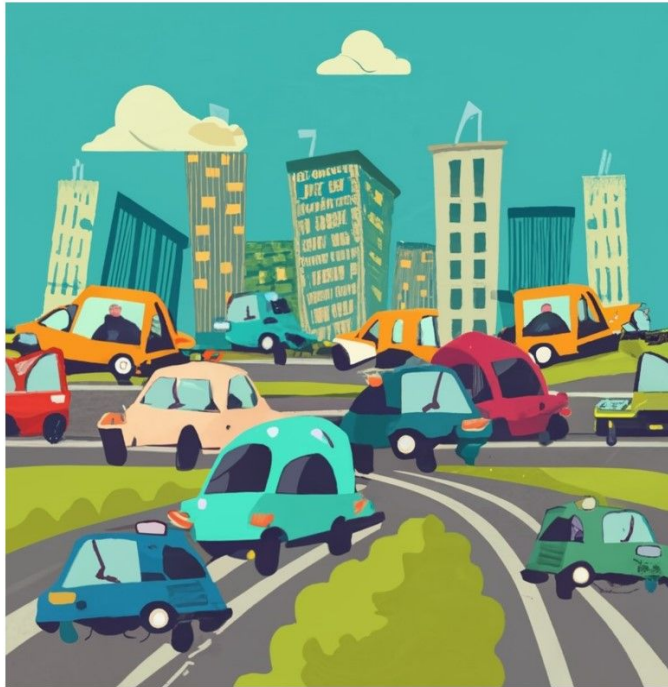


Welcome to Microservice City!



Reference: <https://www.honeycomb.io/microservices>

latency issue!



DTraComp: Distributed Trace Compare

©Maryam Ekhlesi et al.

Presented at ICPE 2024

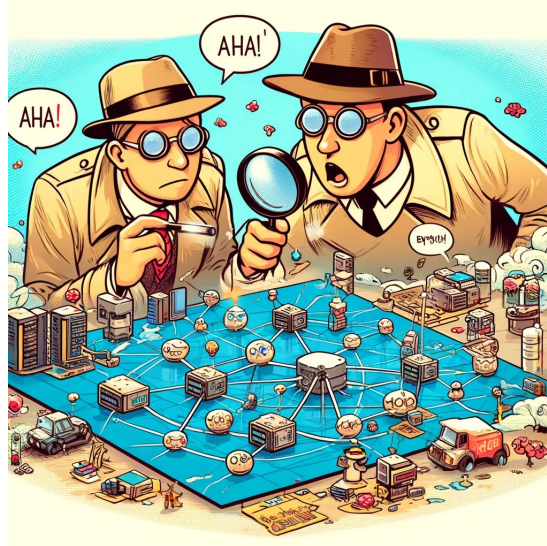
Basic Concepts!

DTraComp: Distributed Trace Compare

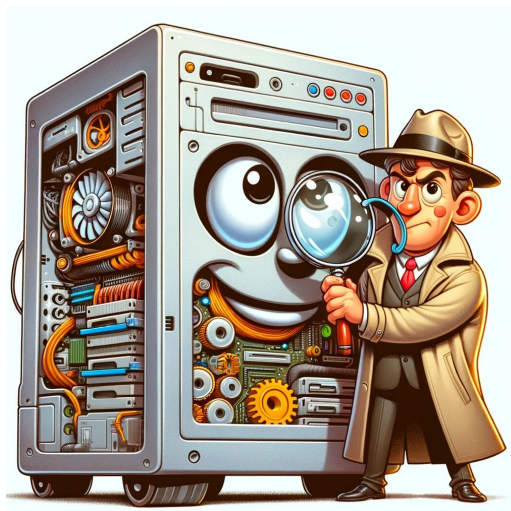
©Maryam Ekhlesi et al.

Presented at ICPE 2024

Question: Have you ever wondered how developers track the behavior of distributed applications in real-time, especially in complex systems?



Question: What technique allows software developers to analyze and monitor the detailed operations of a system, including interactions with the operating system, to enhance performance and diagnose issues?!



Levels of observability

- Debugging
- Logging
- Profiling
- Tracing
 - Low-level Tracing
 - High-level Tracing
- Monitoring



Debugging

- The process of finding errors in source code!
- Performed in the development environment.

```
def calculate_average(numbers):  
    total_sum = sum(numbers)  
    count = len(numbers)  
    average = total_sum / count # Potential division by zero error  
    return average  
  
# List of numbers  
numbers = []  
  
# Calculate the average  
print("The average is:", calculate_average(numbers))
```

 Throw a division by zero **error!!**

Logging

- Recording information during the execution time.
- Levels of logging
 - DEBUG
 - INFO
 - WARNING
 - ERROR
 - FATAL

```
import logging

# Configure logging
logging.basicConfig(level=logging.DEBUG, format='%(asctime)s - %(levelname)s -

def process_transaction(transaction_id, amount):
    logging.info(f"Starting transaction for ID: {transaction_id} with amount:
    try:
        # Simulate a transaction processing
        if amount <= 0:
            raise ValueError("Amount must be positive")
        # Simulate successful transaction processing
        logging.info(f"Transaction {transaction_id} processed successfully")
    except ValueError as e:
        logging.error(f"Error processing transaction {transaction_id}: {e}")

# Example of processing transactions
process_transaction('TX1001', 500)
process_transaction('TX1002', -150)
```

Profiling

- Statistical summary of observed events.
- Where is performance lost?
- Won't report the reason!

```
import cProfile
import random
import time

def process_data_fast():
    time.sleep(1) # Simulates processing time
    return sum(random.sample(range(10000), 1000))

def process_data_medium():
    time.sleep(2) # Simulates more processing time
    return sum(random.sample(range(10000), 5000))

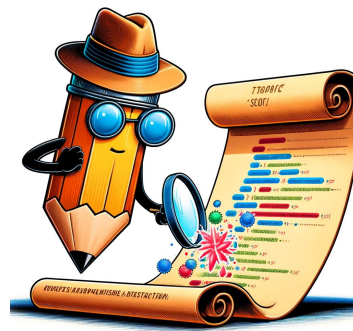
def process_data_slow():
    time.sleep(3) # Simulates even more processing time
    data = [random.randint(1, 100) for _ in range(10000)]
    return sum(data)

def main():
    fast_result = process_data_fast()
    medium_result = process_data_medium()
    slow_result = process_data_slow()
    print(f"Fast: {fast_result}, Medium: {medium_result}, Slow: {slow_result}")

if __name__ == "__main__":
    profiler = cProfile.Profile()
    profiler.enable()
    main()
    profiler.disable()
    profiler.print_stats(sort='time')
```

Low-level Tracing!

- Printf!
- Every location in the code that we want to trace is referred to as a **tracepoint**. The process of inserting tracepoints into the code is known as **instrumentation** and each time a tracepoint is executed, it is generated an **event**.
- Instrumentation
 - Static instrumentation
 - Dynamic instrumentation



Low-level Tracing!

```
#include <stdio.h>

#define MAX 5

int calculate_sum(int* array) {
    int total = 0;

    printf("Starting calculation");           // <-- tracepoint
    for (int i = 0; i < MAX; i++) {
        total += array[i];
        printf("Adding index %d, total now %d", i, total); // <-- tracepoint
    }
    printf("Final total: %d", total);        // <-- tracepoint

    return total;
}

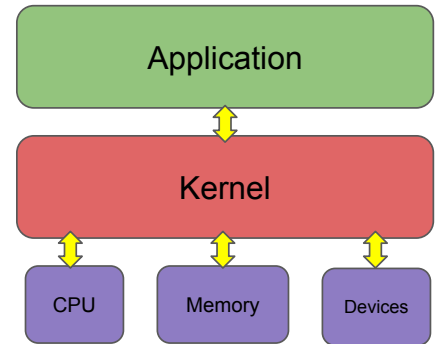
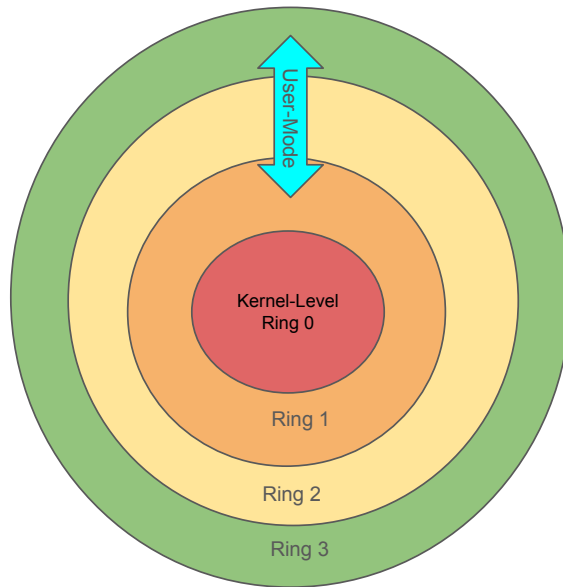
int main() {
    int numbers[MAX] = {1, 2, 3, 4, 5};
    calculate_sum(numbers);
    return 0;
}
```

```
$ ./myprog
Starting calculation           // <-- event
Adding index 0, total now 1   // <-- event
Adding index 1, total now 3   // <-- event
Adding index 2, total now 6   // <-- event
Adding index 3, total now 10  // <-- event
Adding index 4, total now 15  // <-- event
Final total: 15              // <-- event
```


Two Levels of low-level tracing

User-space Tracing

Kernel-Level Tracing

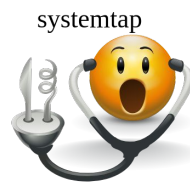


Tracing VS Logging

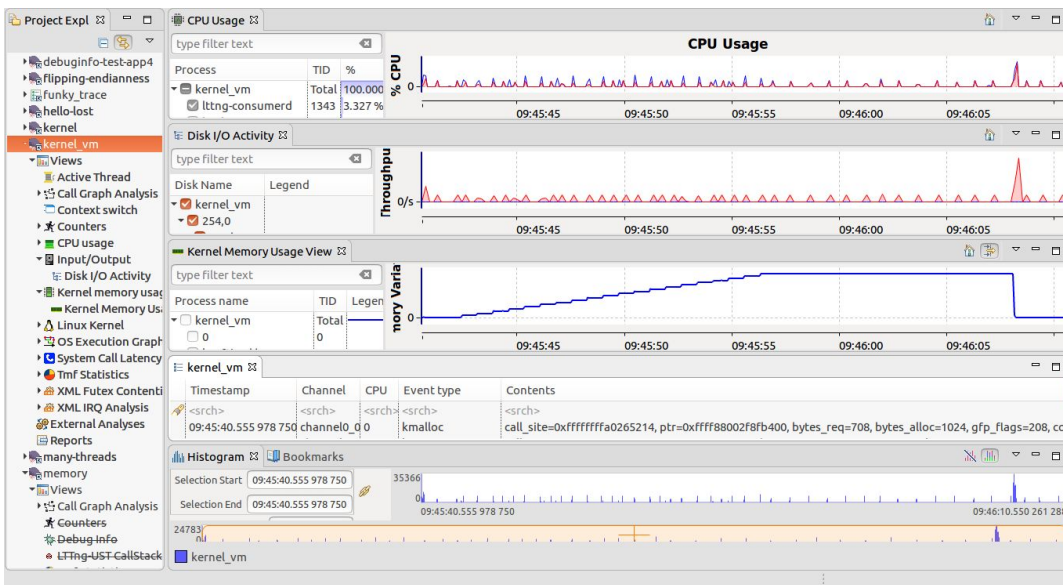
- Tracers are designed to record much lower-level events that occur much more frequently than log messages.
- Logging is more appropriate for a very high-level analysis of less frequent events: user accesses, exceptional conditions (errors and warnings, for example), database transactions, instant messaging communications, and such. Simply put, logging is one of the many use cases that can be satisfied with tracing.

Competing software tracers (low-level tracers)

- LTTng
- Dtrace4linux
- eBPF
- Ftrace
- Perf
- Strace
- Sysdig
- SystemTap



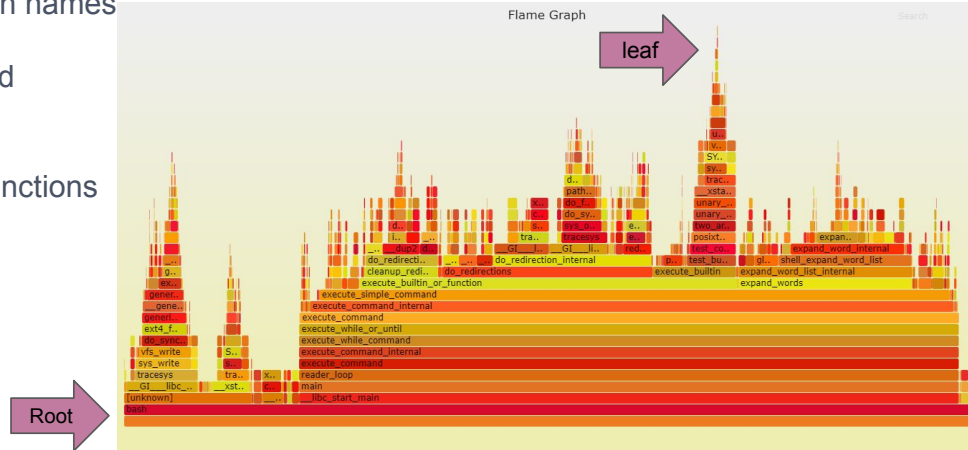
Low-level Trace Visualization Tool



Visualization Approaches

Flame Graph

- The horizontal axis of the call stack collection represents the function names
- Alphabetical arrangement of function names from left to right
- Combining identical functions placed side by side horizontally
- The width of each box shows the appearance frequency of that functions
- Random colors



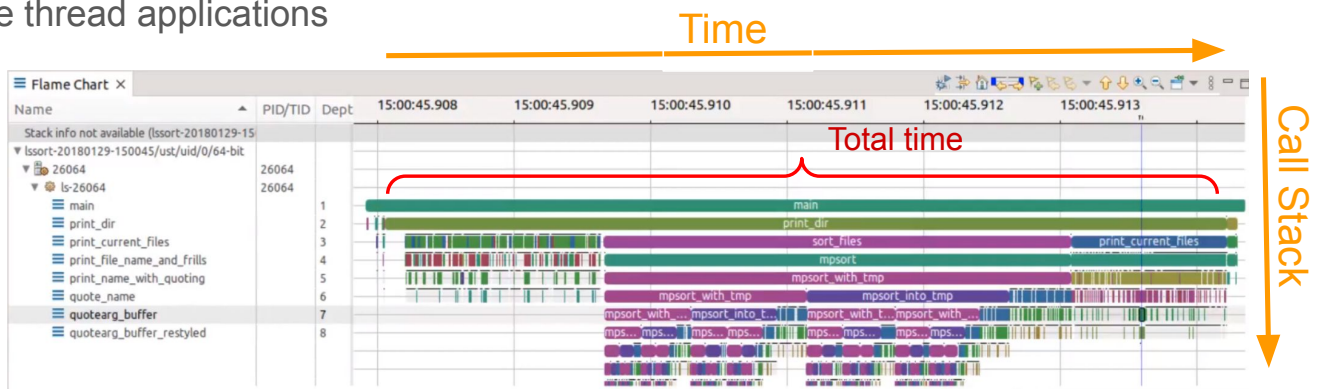
Reference: <https://github.com/brendangregg/FlameGraph>

Visualization Approaches

Flame Chart

- Time on the horizontal axis
- Call stack on the vertical axis
- Single thread applications

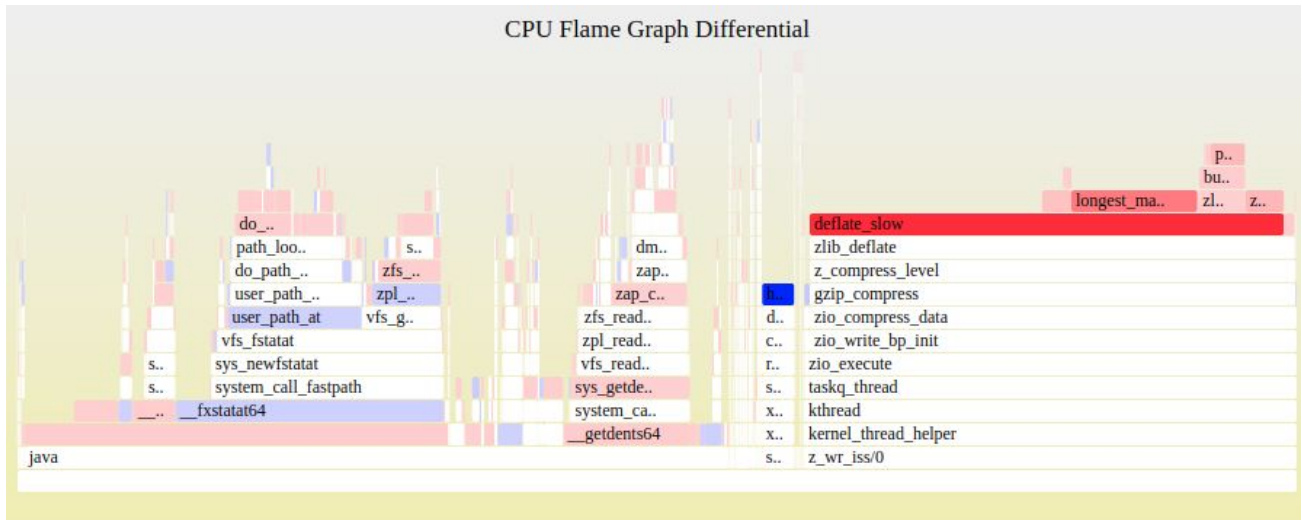
The self time of a function is indicated by the difference in width between its frame and the frames directly beneath it



Reference: <https://archive.eclipse.org/tracecompass/doc/stable/org.eclipse.tracecompass.doc.user/LTTng-UST-Analyses.html>

Visualization Approaches

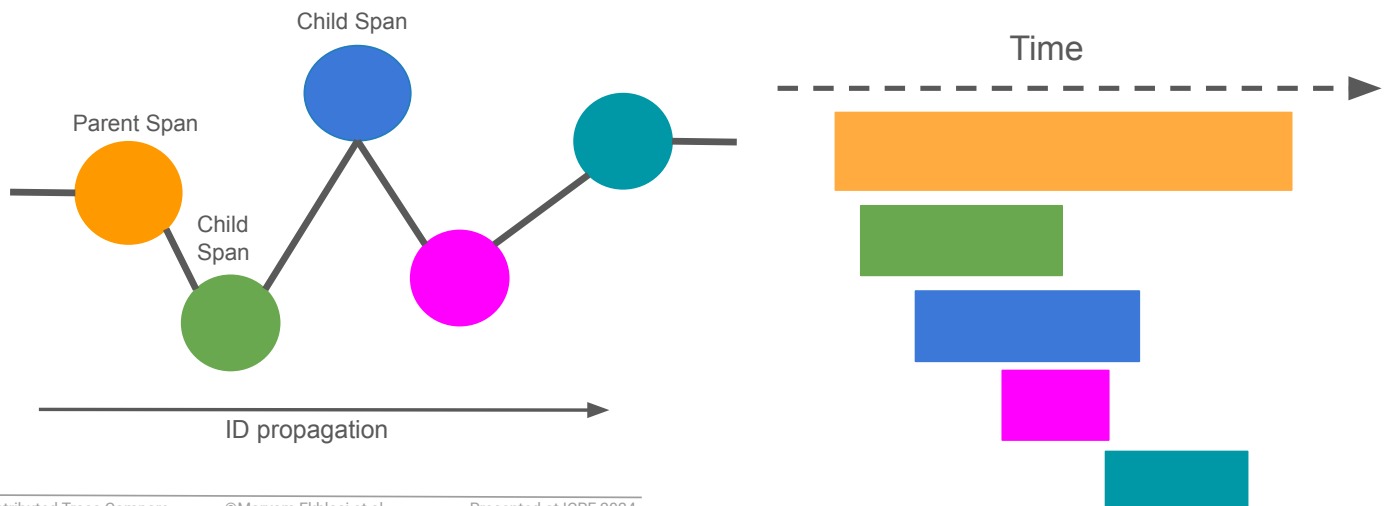
Differential Flame Graph



Reference: <https://www.brendangregg.com/blog/2014-11-09/differential-flame-graphs.html>

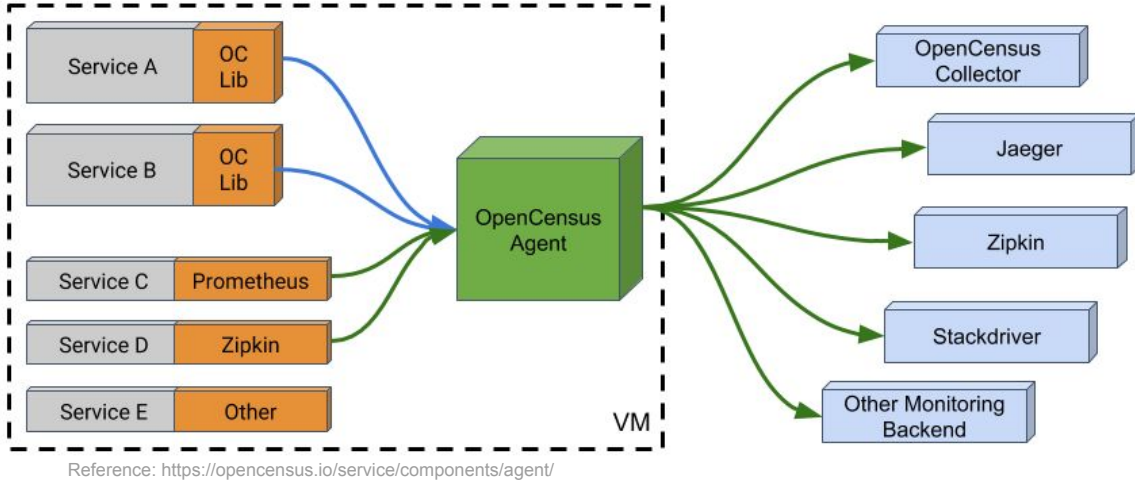
Distributed Tracing (High-level Tracing)

- End-to-end tracing!
- Technique for monitoring application requests as they move from frontend devices through backend services, databases, and any intermediary services.



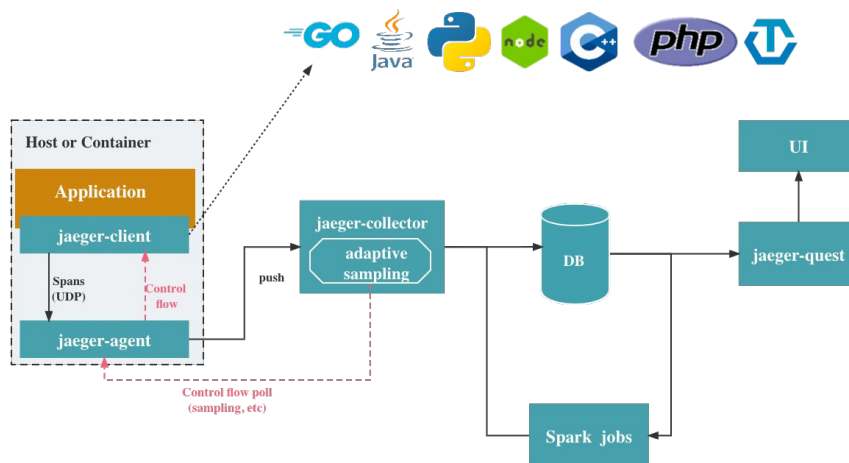
Distributed End-to-end Tracing Standards

1. OpenCensus



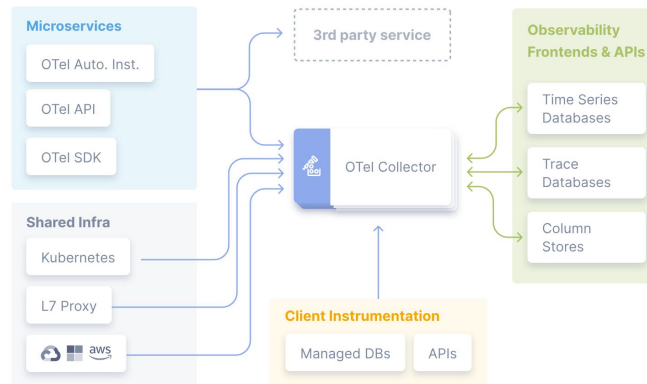
Distributed End-to-end Tracing Standards

2. OpenTracing



Distributed End-to-end Tracing Standards

3. OpenTelemetry

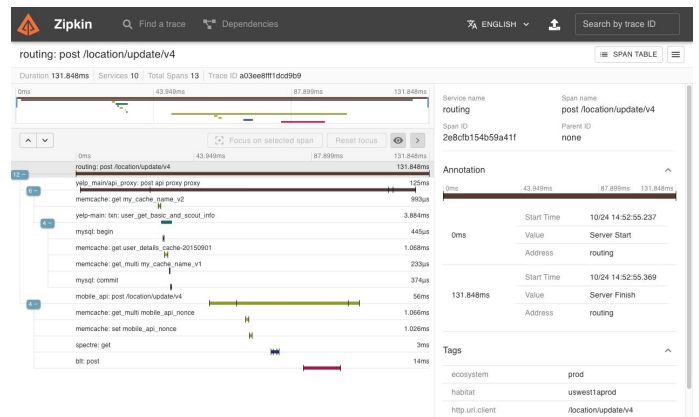
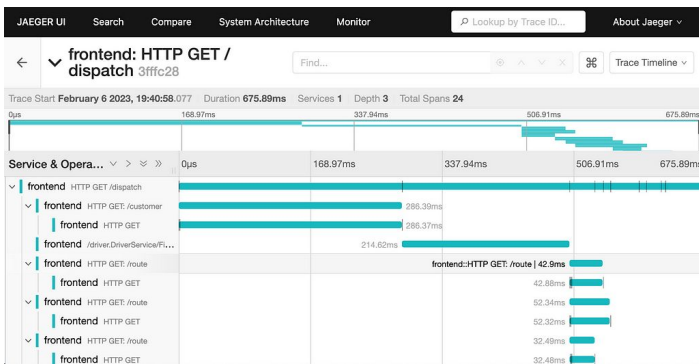


Reference: <https://opentelemetry.io/docs/>

High-level Trace Visualization Tools



- Visualizing distributed requests
- Distributed context propagation
- Span's parent-child relationships



Power of Tracing!

- Over Profiling
- Over Debugging
- Over Logging
- Over Monitoring

Profiling	Tracing
Execution	Execution
	Latency Spike Detected

Software Performance!



Deadlines!



Fast!



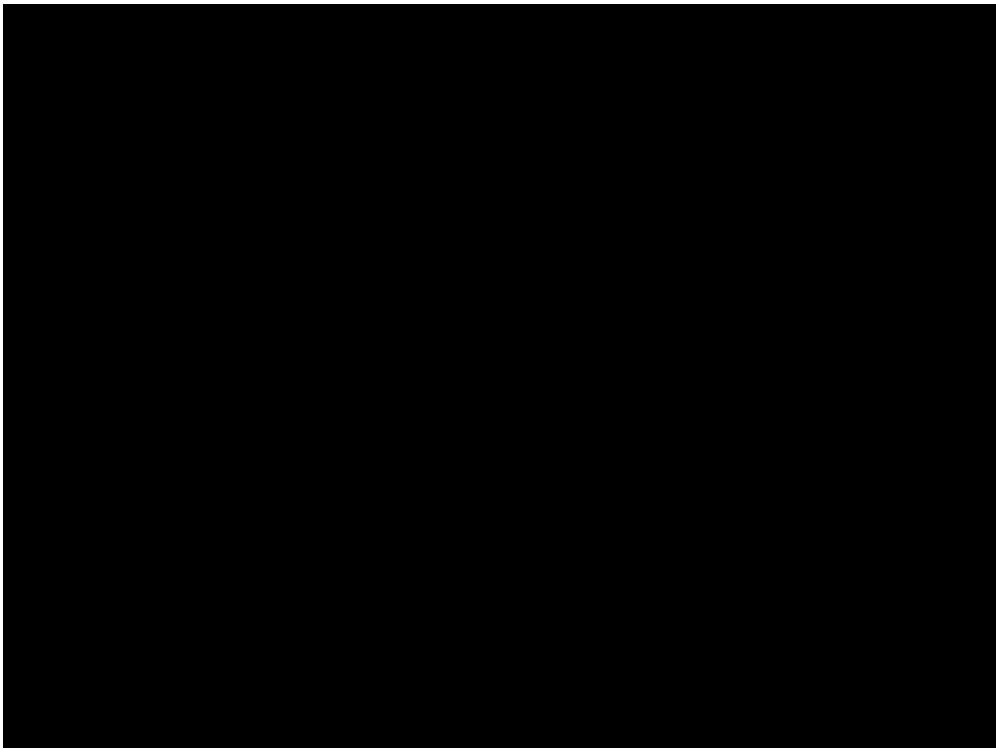
Limited Machine Resources!

What is Software Performance

- Software Aspects
 - **Code efficiency**
 - Caching strategy
 - Architecture and Design
 - Concurrency and Parallelism
- Hardware Aspects
 - Processor Speed Core Count
 - Memory (RAM)
 - Storage (HDD/SSD)
 - Network Speed
- Other Factors



Motivation! Code Efficiency!



Executions Differences!



✓ TiDB: server.dispatch 6e96e6d
808µs

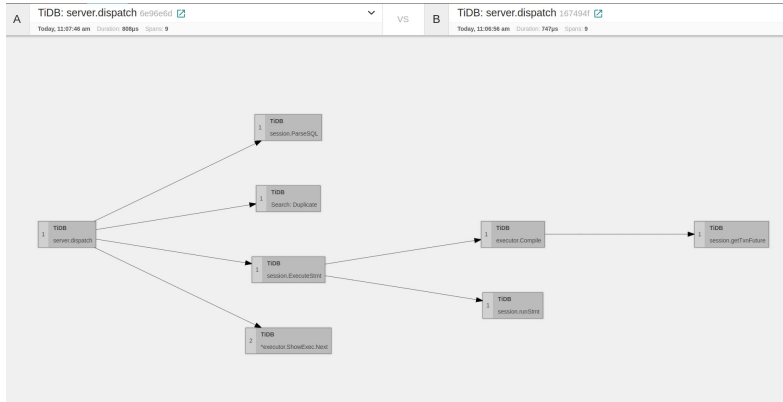
9 Spans
TiDB (9)

Today 11:07:46 am
13 minutes ago

✓ TiDB: server.dispatch 167494f
747µs

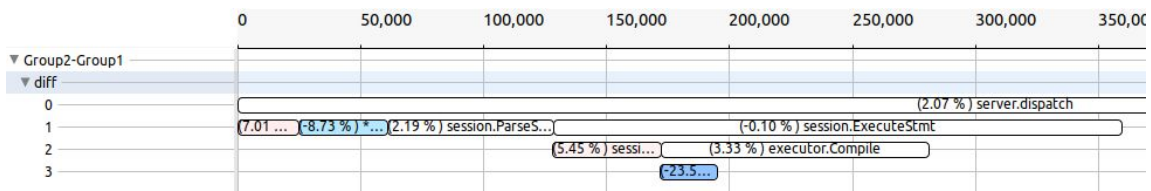
9 Spans
TiDB (9)

Today 11:06:56 am
14 minutes ago



DTraComp: Distributed Trace Compare ©Maryam Ekhlesi et al. Presented at ICPE 2024

Detective DTraComp Enter!



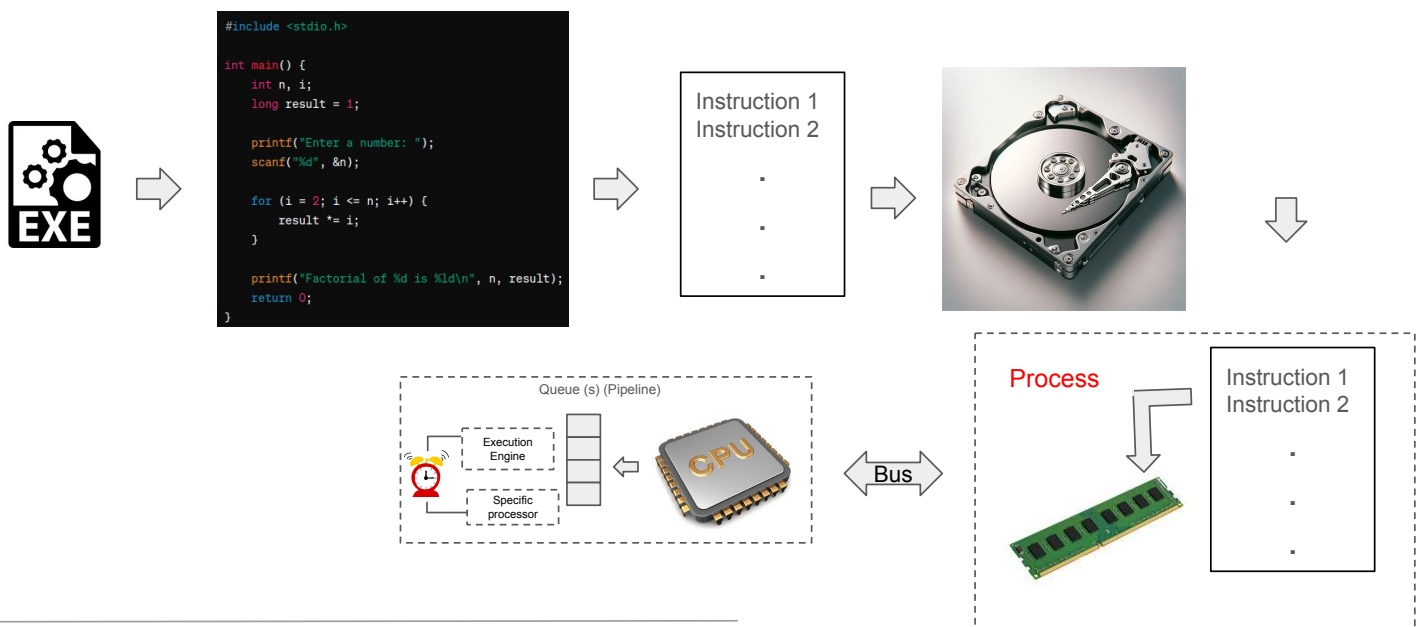
Start Time	End Time	Duration	SpanId	OperationName	TraceID	ThreadID	RUN	RUNSYSTEMCALL	INTERRUPTED	WAITCPU	WAITUNKNOWN	WAITFORK	UNKNOWN
11:06:56.933 952 078	11:06:56.933 980 155	28,077	4442450064020525663	*executor.ShowExec.Next	1618081973820362342	1805943	17108	12915	0	0	0	0	0
11:06:56.934 022 060	11:06:56.934 042 663	20,603	6023047575761396435	*executor.ShowExec.Next	1618081973820362342	1805943	7188	15709	0	0	0	0	0
11:07:46.711 230 210	11:07:46.711 249 416	19,206	2871436018761349651	*executor.ShowExec.Next	7968810417120810813	1813228	10054	11450	0	0	0	0	0
11:07:46.711 292 718	11:07:46.711 303 334	10,616	1574938190389480438	*executor.ShowExec.Next	7968810417120810813	1813228	4326	8727	0	0	0	0	0
11:06:56.933 723 277	11:06:56.933 843 056	119,779	7456353795818001242	executor.Compile	1618081973820362342	1805943	68216	53966	0	0	0	0	0
11:07:46.711 033 606	11:07:46.711 132 920	99,314	2159668538522373865	executor.Compile	7968810417120810813	1813228	52290	49358	0	0	0	0	0
11:06:56.933 567 600	11:06:56.933 592 324	24,724	3369336188777566051	Search: Duplicate	1618081973820362342	1805943	10539	16130	0	0	0	0	0
11:07:46.710 746 696	11:07:46.710 772 468	25,772	8310292673920133684	Search: Duplicate	7968810417120810813	1813228	10960	16687	0	0	0	0	0

We need more detail informations to answer this question!

What is happening inside the system?!

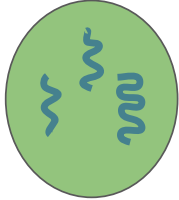
- How a program runs?
- What are the processes and threads?
- What are kernel and user space levels?
- What are process/ thread states?

What is a process?



What is a process?

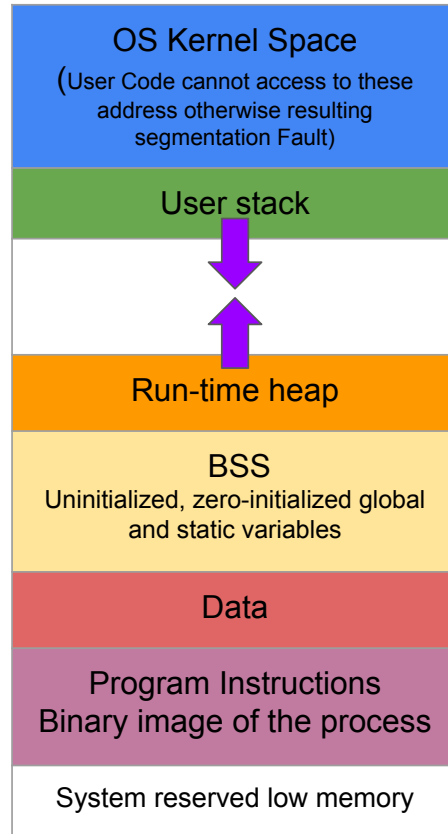
Program in execution!!



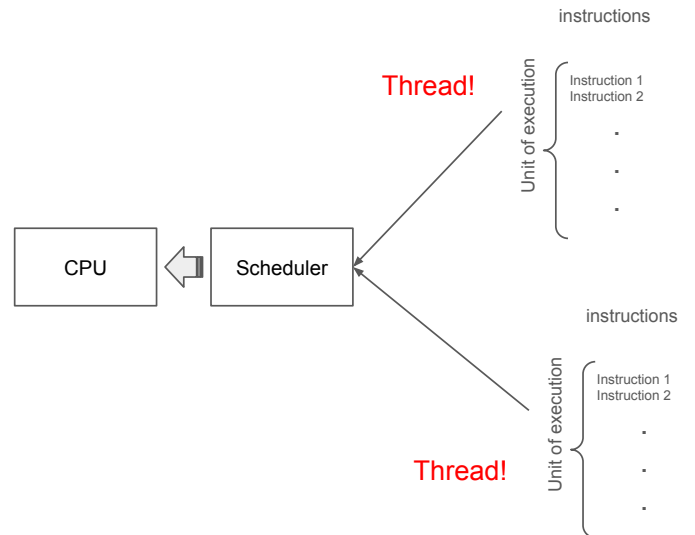
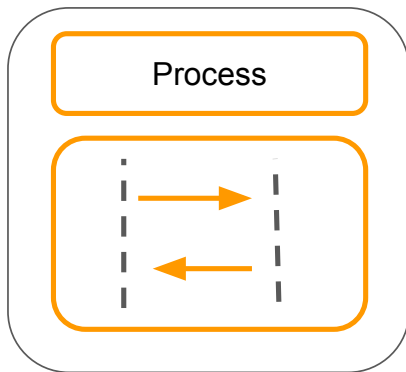
Local variables in functions,
function call management

Objects, dynamic memory

Fixed size!



What is a thread?



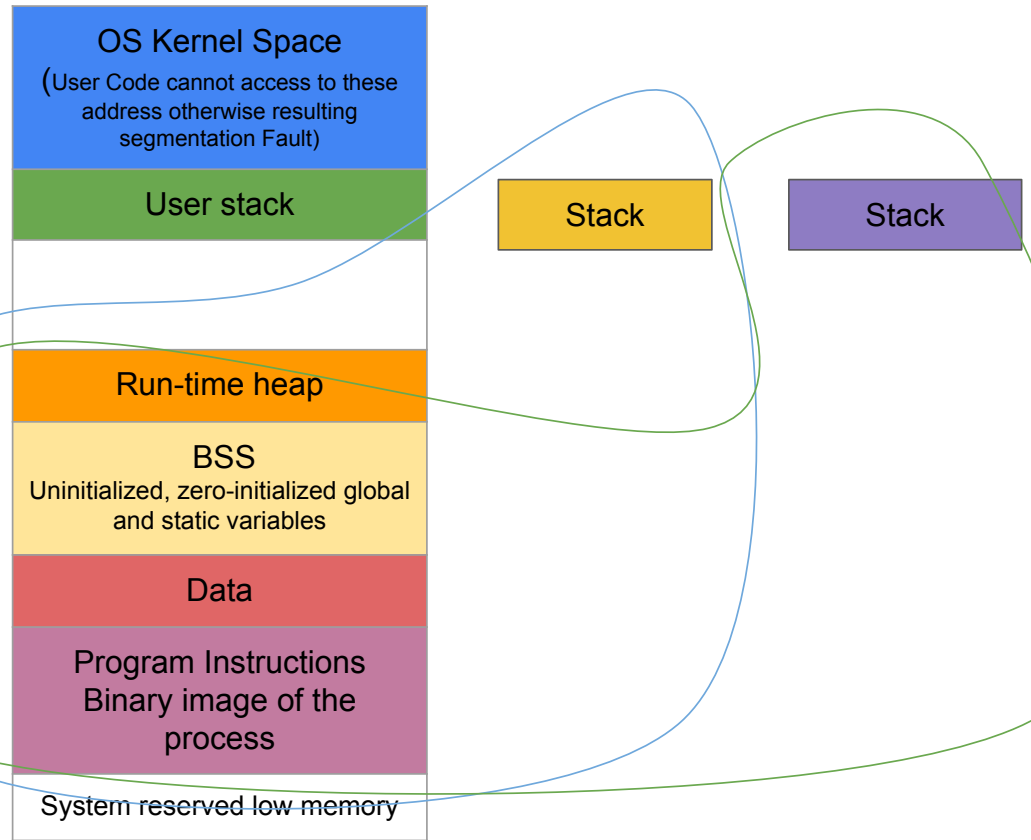
What is a thread?

A Thread contains

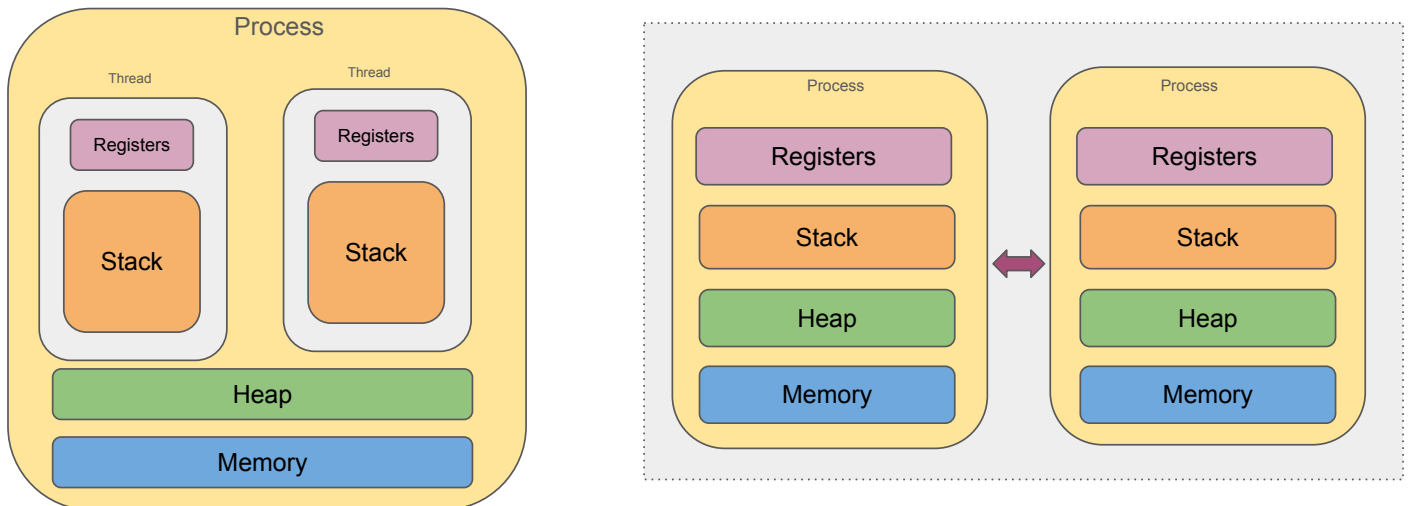
- A thread ID
- A program counter
- A register set
- A stack

Benefits of threads

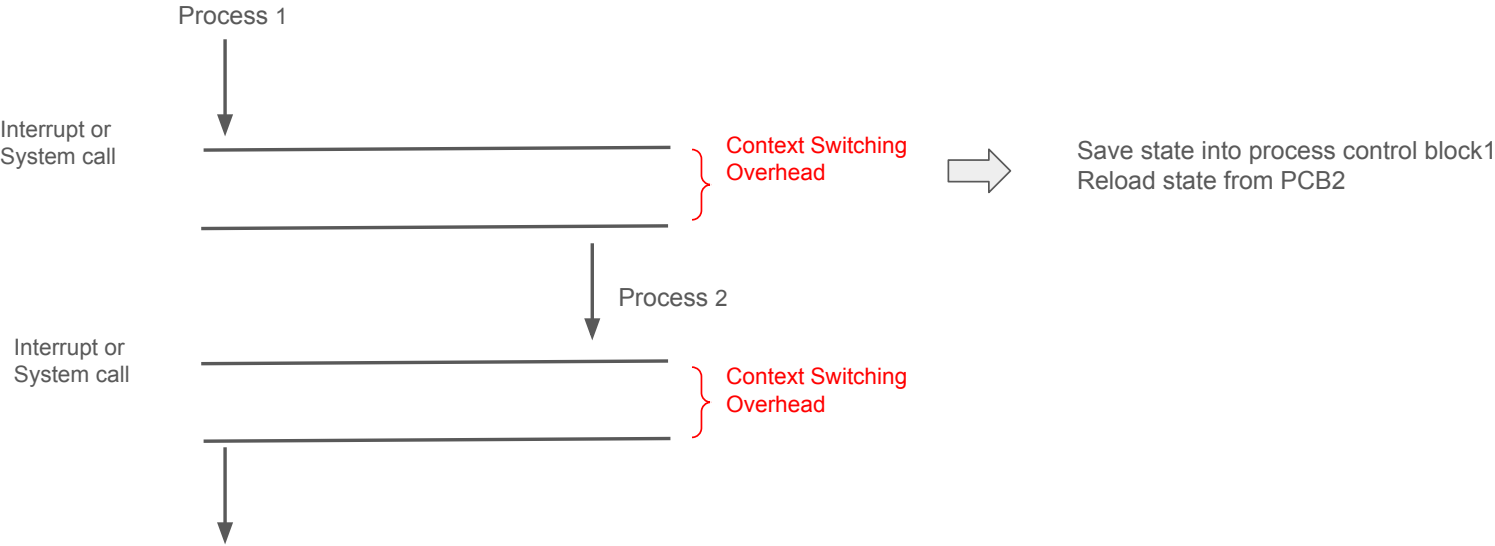
- Each process can do more than one task at a time
- Threads are lightweight however processes are heavyweight



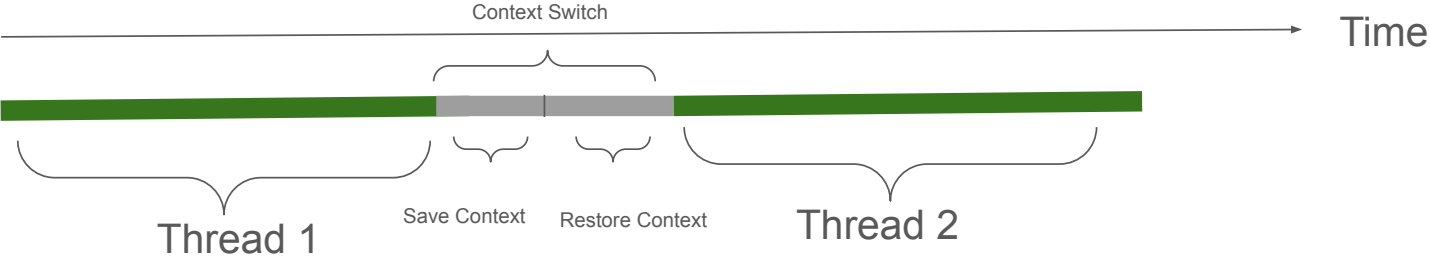
Differences between threads and Processes



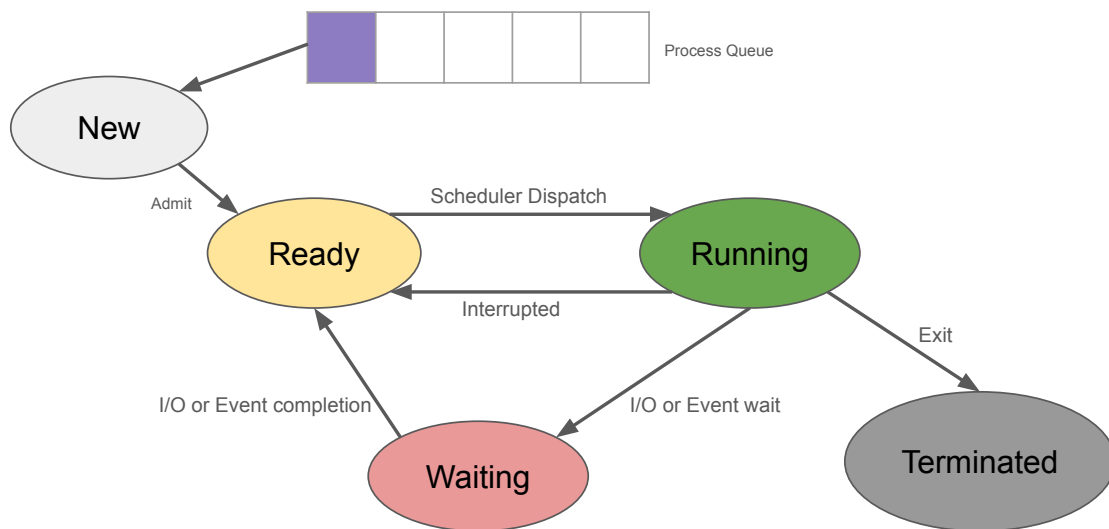
Context Switch in processes



Context Switch in threads

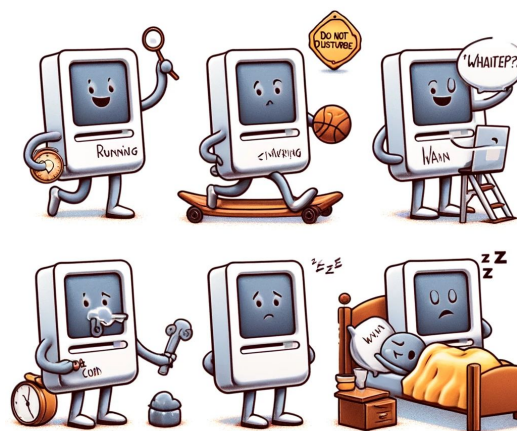


We have to know the process states

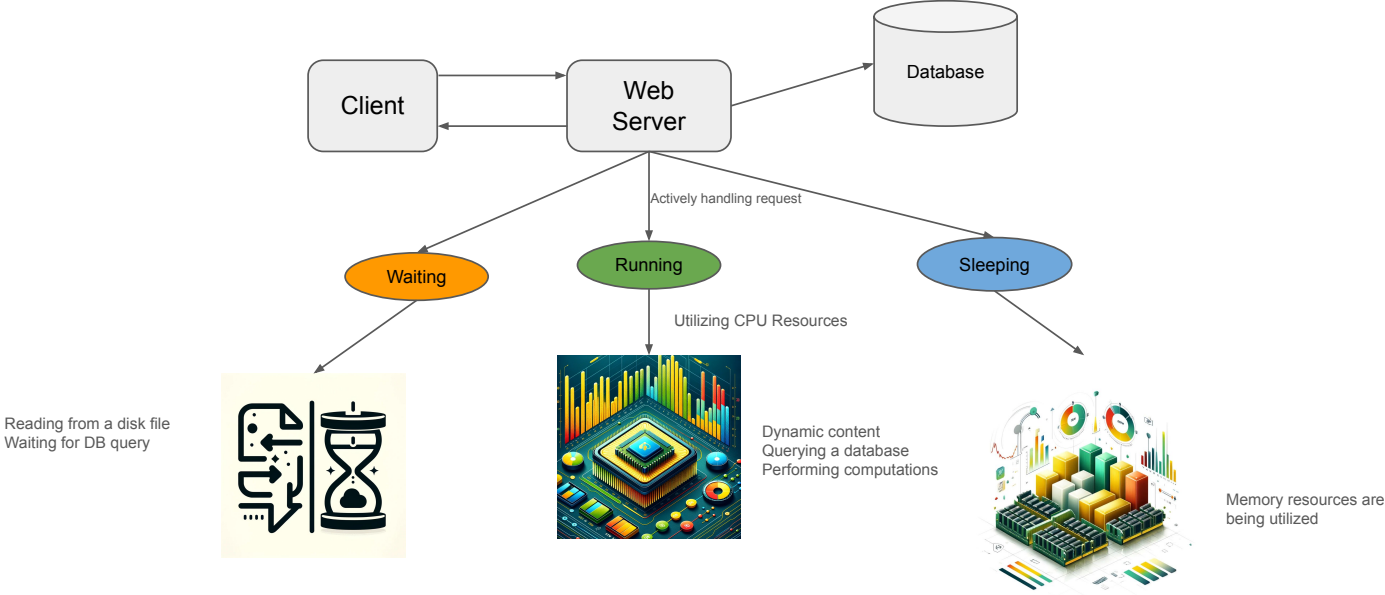


What the state of a process can tell us

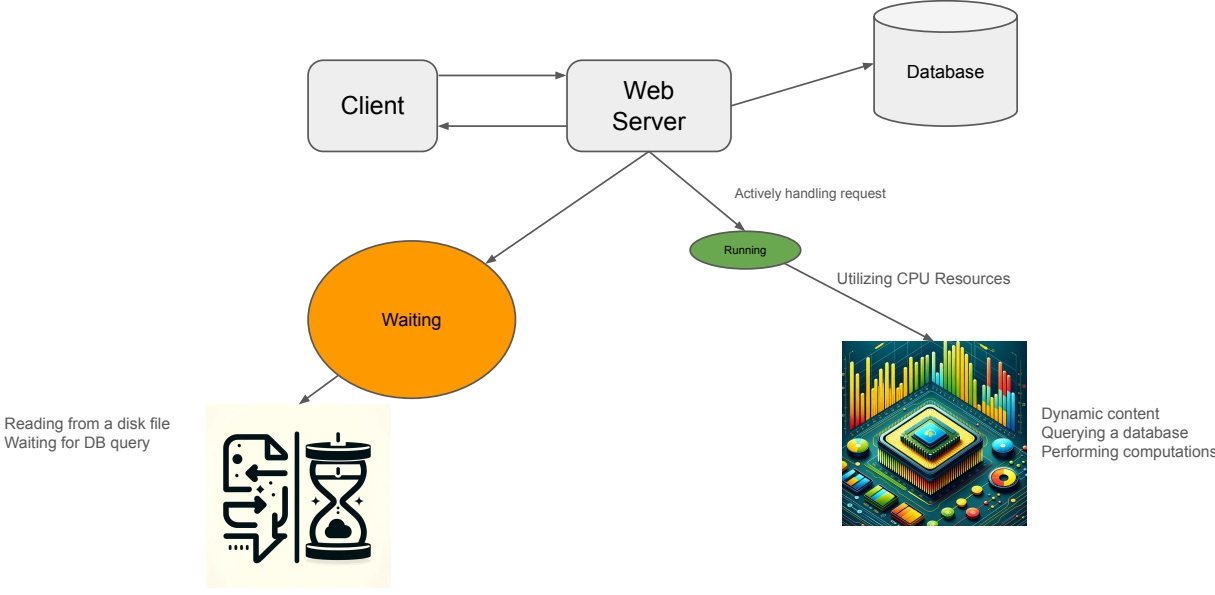
- Process Status
- Resource Utilization
- Prioritization
- I/O Status
- Inter-process Communication (IPC)
- Concurrency and Parallelism
- Life Cycle Management



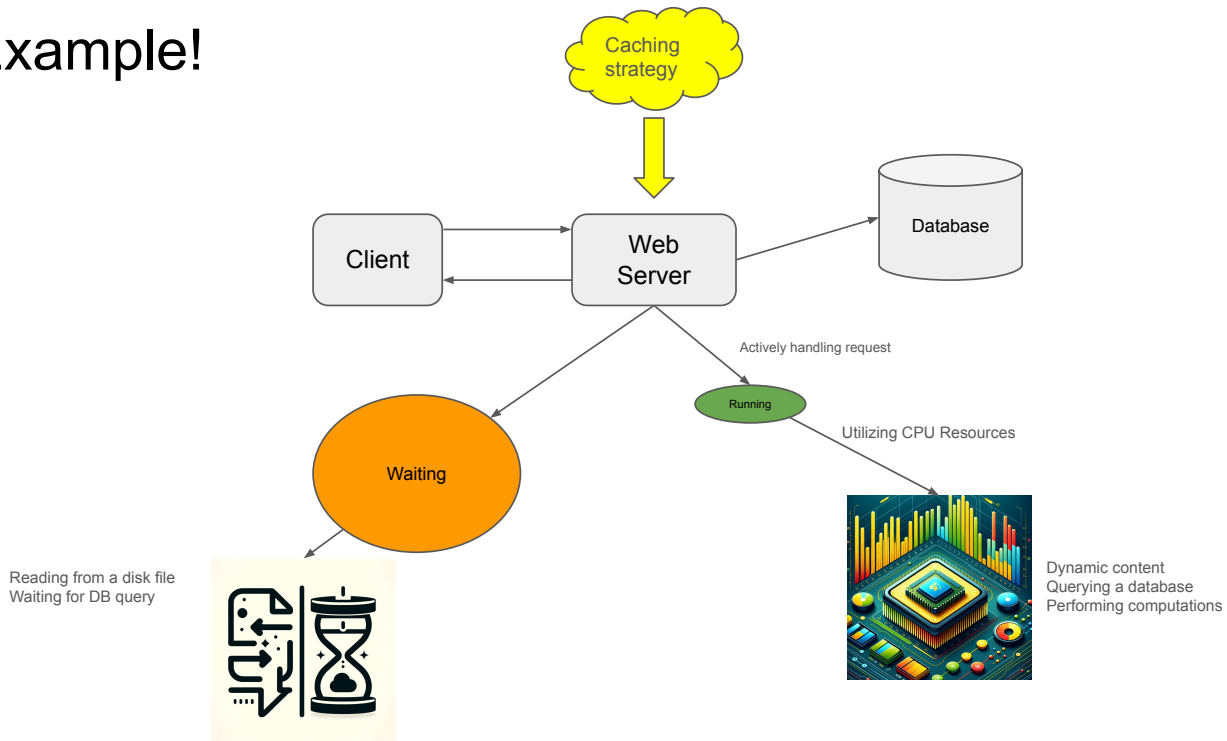
Example!



Example!

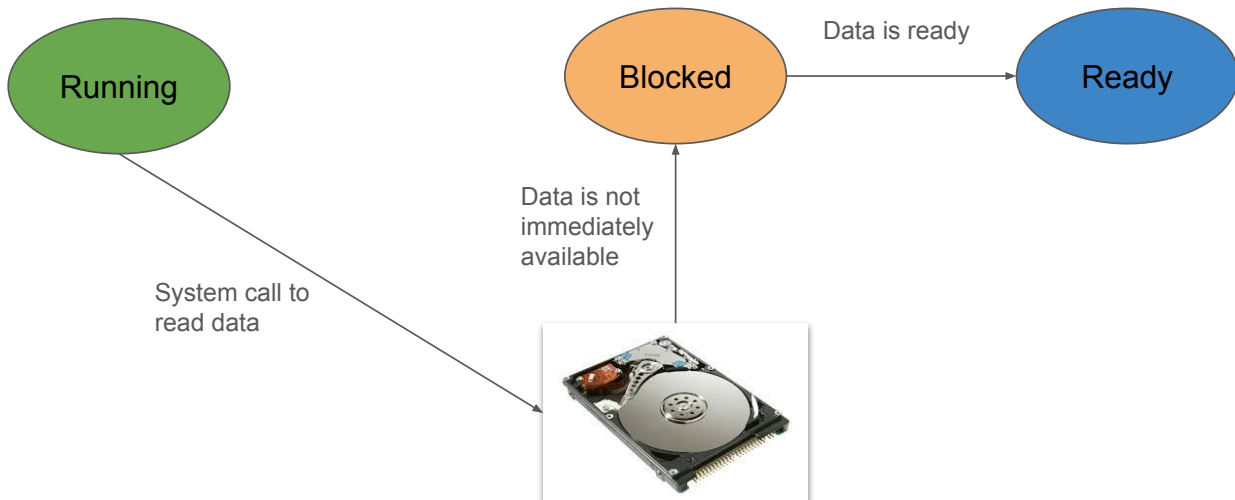


Example!

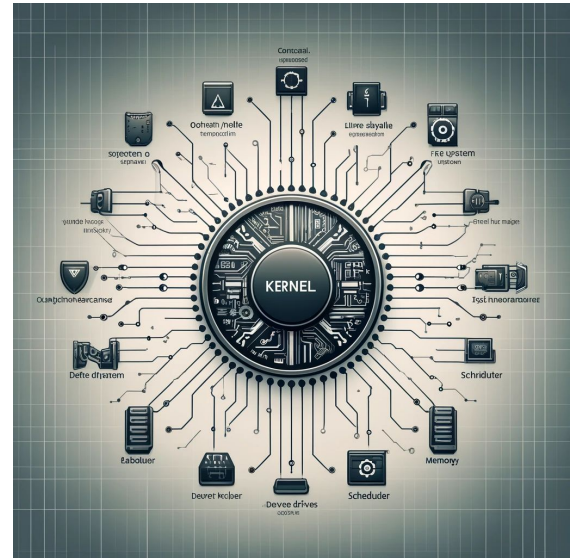
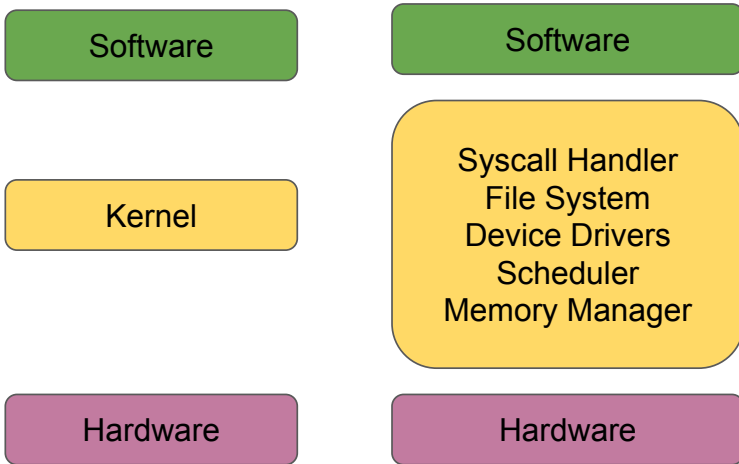


How we can know why the process is in the waiting state?

By system calls



What is Kernel?



What are system calls?

- Process Control
- File Manipulation
- Device Management
- Information Maintenance
- Communication



We have all the required information let's have an example!

What is the time complexity?

Complexity: Using Nested Loops

Complexity: Using a Hash Table

```
// Function to check for duplicates in an array using nested loops
bool hasDuplicates(int arr[], int size) {
    for (int i = 0; i < size - 1; i++) {
        for (int j = i + 1; j < size; j++) {
            if (arr[i] == arr[j]) {
                return true; // Duplicate found
            }
        }
    }
    return false; // No duplicates
}
```

```
// Function to check for duplicates in an array using a simplistic hash table
bool hasDuplicatesOptimized(int arr[], int size) {
    bool hashTable[MAX_VALUE] = { false };

    for (int i = 0; i < size; i++) {
        if (arr[i] < 0 || arr[i] >= MAX_VALUE) {
            printf("Value out of range.\n");
            return false; // Value out of the assumed range
        }
        if (hashTable[arr[i]]) {
            return true; // Duplicate found
        }
        hashTable[arr[i]] = true;
    }
    return false; // No duplicates
}
```

Actual Use case - C Language

$O(n^2)$ Complexity: Using Nested Loops

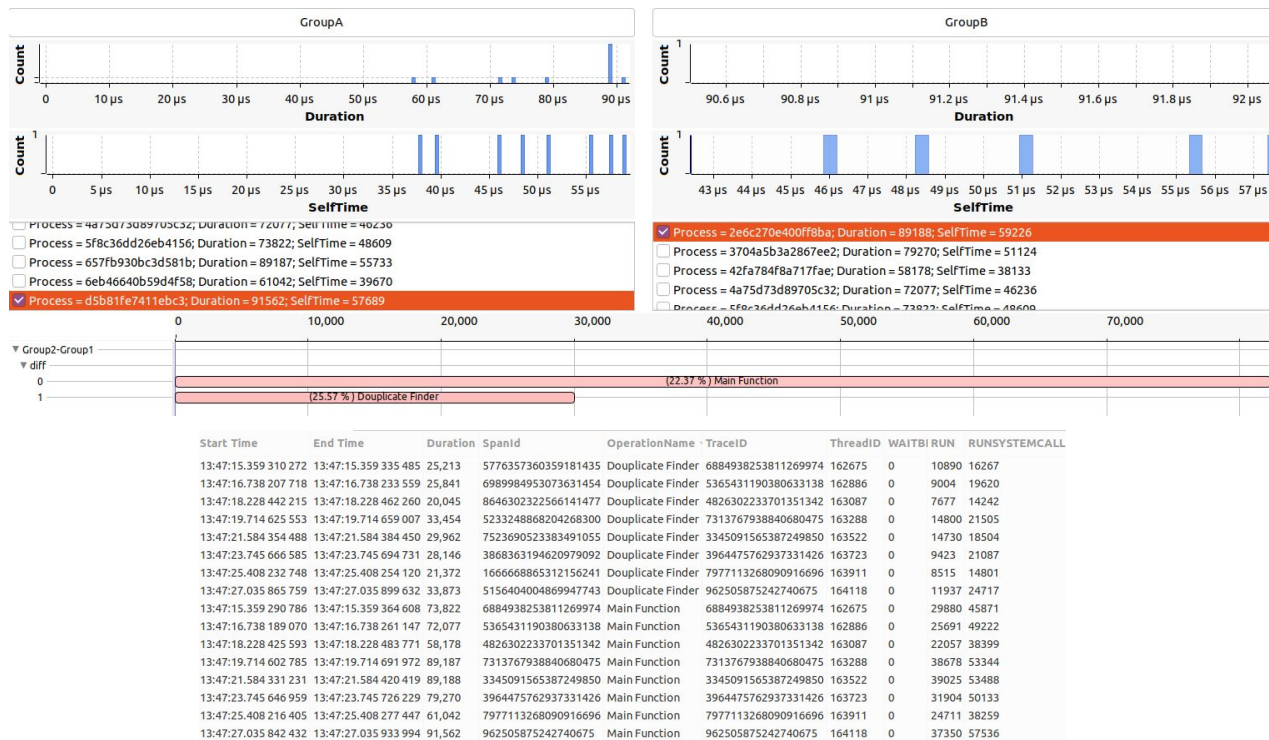
$O(n)$ Complexity: Using a Hash Table

```
// Function to check for duplicates in an array using nested loops
// Time Complexity: O(n^2)
bool hasDuplicates(int arr[], int size) {
    for (int i = 0; i < size - 1; i++) {
        for (int j = i + 1; j < size; j++) {
            if (arr[i] == arr[j]) {
                return true; // Duplicate found
            }
        }
    }
    return false; // No duplicates
}
```

```
// Function to check for duplicates in an array using a simplistic hash table
// Time Complexity: O(n), with the caveat of known and limited range of input values
bool hasDuplicatesOptimized(int arr[], int size) {
    bool hashTable[MAX_VALUE] = { false };

    for (int i = 0; i < size; i++) {
        if (arr[i] < 0 || arr[i] >= MAX_VALUE) {
            printf("Value out of range.\n");
            return false; // Value out of the assumed range
        }
        if (hashTable[arr[i]]) {
            return true; // Duplicate found
        }
        hashTable[arr[i]] = true;
    }
    return false; // No duplicates
}
```

Differential Flame graph



What is Software Performance

- Software Aspects
 - Code efficiency
 - **Caching strategy**
 - Architecture and Design
 - Concurrency and Parallelism
- Hardware Aspects
 - Processor Speed Core Count
 - Memory (RAM)
 - Storage (HDD/SSD)
 - Network Speed
- Other Factors

Caching Strategy

No-caching

```
import (  
    "database/sql"  
    "fmt"  
    "time"  
)  
    _ "github.com/mattn/go-sqlite3"  
  
// User struct to hold user details  
type User struct {  
    ID int  
    Name string  
}  
  
// Global cache variable  
var userCache = make(map[int]User)  
  
// Non-caching version  
func getUserDetailsWithoutCache(userID int) (*User, error) {  
    db, err := sql.Open("sqlite3", "users.db")  
    if err != nil {  
        return nil, err  
    }  
    defer db.Close()  
  
    var user User  
    err = db.QueryRow("SELECT id, name FROM users WHERE id = ?", userID).Scan(&user.ID, &user.Name)  
    if err != nil {  
        return nil, err  
    }  
  
    return &user, nil  
}
```

Execution time without cache: 0.05859804153442383 seconds
Execution time with cache: 0.0003237724304199219 seconds

Caching

```
import (  
    "database/sql"  
    "fmt"  
    "time"  
)  
    _ "github.com/mattn/go-sqlite3"  
  
// User struct to hold user details  
type User struct {  
    ID int  
    Name string  
}  
  
// Global cache variable  
var userCache = make(map[int]User)  
  
// Caching version  
func getUserDetailsWithCache(userID int) (*User, error) {  
    if user, found := userCache[userID]; found {  
        return &user, nil  
    }  
  
    user, err := getUserDetailsWithoutCache(userID)  
    if err != nil {  
        return nil, err  
    }  
  
    userCache[userID] = *user  
    return user, nil  
}
```

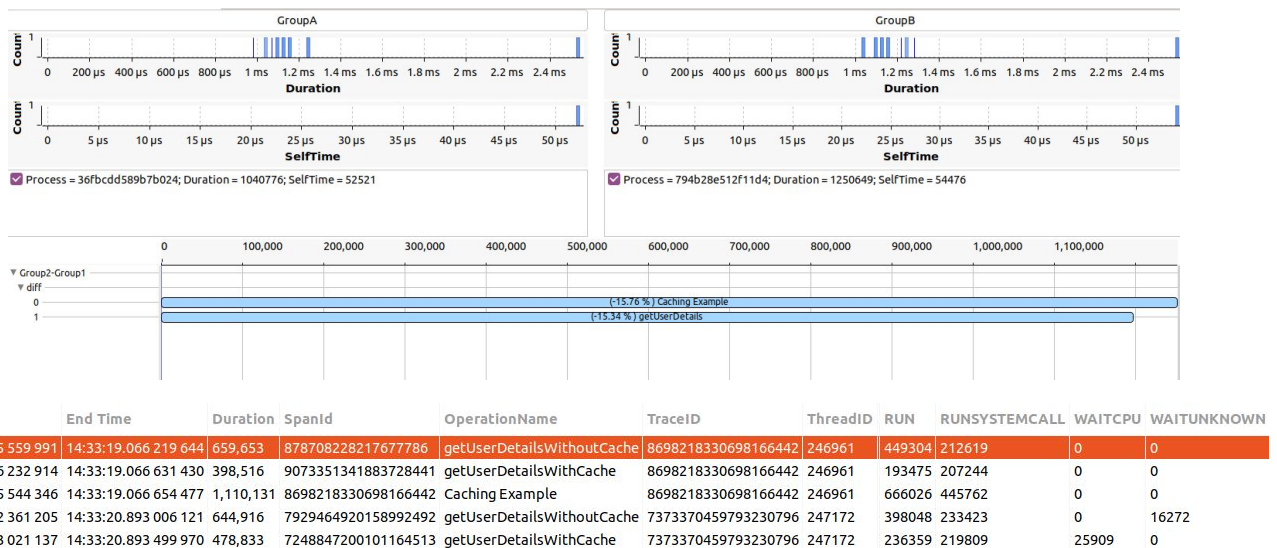
Expectation: Increase in blocked time and wait time since we need to fetch data from database
More context switch and cpu execution time for handling connection

What do you expect?

- Without Caching: Expect more frequent and longer periods in the I/O wait and possibly blocked states
- With Caching: The frequency and duration of I/O wait and blocked states are significantly reduced

Expectations are Confirmed!

- Without Caching: Expect more frequent and longer periods in the I/O wait and possibly blocked states
- With Caching: The frequency and duration of I/O wait and blocked states are significantly reduced

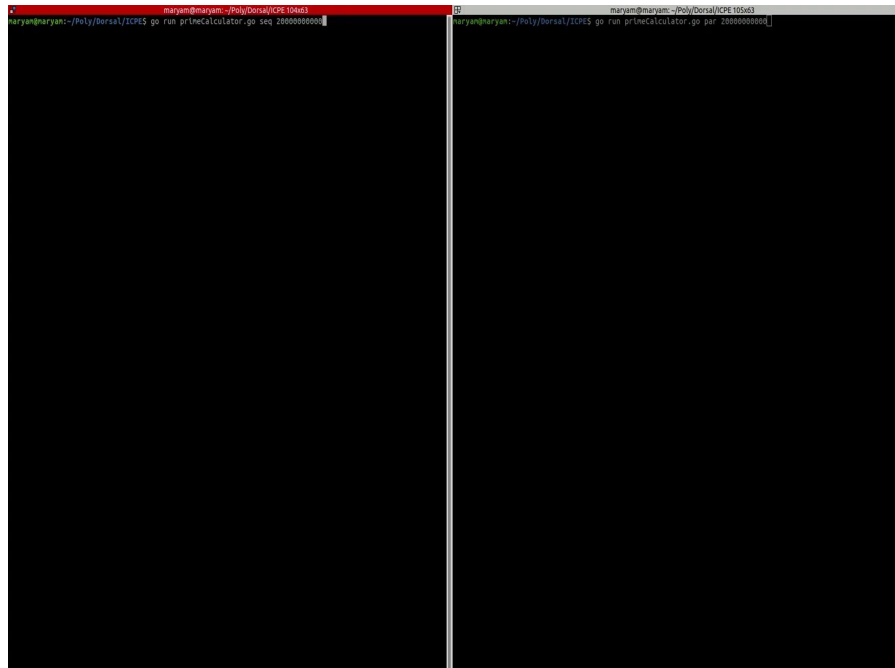


What is Software Performance

- Software Aspects
 - Code efficiency
 - Caching strategy
 - Architecture and Design
 - **Concurrency and Parallelism**
- Hardware Aspects
 - Processor Speed Core Count
 - Memory (RAM)
 - Storage (HDD/SSD)
 - Network Speed
- Other Factors

Concurrency and Parallelism

CPU Bound Example:



Parallel and sequential

```
// isPrime checks if a number is a prime
func isPrime(number int) bool {
    if number <= 1 {
        return false
    }
    for i := 2; i <= int(math.Floor(math.Sqrt(float64(number)))); i++ {
        if number%i == 0 {
            return false
        }
    }
    return true
}

// calculatePrimesSeq calculates the first n primes sequentially
func calculatePrimesSeq(upperLimit int) []int {
    primes := []int{}
    for num := 2; num <= upperLimit; num++ {
        if isPrime(num) {
            primes = append(primes, num)
            fmt.Printf("Found Sequential Prime: %d\n", num) // Print found primes
        }
    }
    return primes
}
```

```
// isPrime checks if a number is a prime
func isPrime(number int) bool {
    if number <= 1 {
        return false
    }
    for i := 2; i <= int(math.Floor(math.Sqrt(float64(number)))); i++ {
        if number%i == 0 {
            return false
        }
    }
    return true
}

// calculatePrimesParallel calculates the first n primes in parallel
func calculatePrimesParallel(upperLimit int) []int {
    var primes []int
    var mu sync.Mutex
    var wg sync.WaitGroup

    // Use a channel to generate numbers to be checked for primality
    numbers := make(chan int, 100) // Buffered channel

    // Start a fixed number of worker goroutines to check numbers for primality in parallel
    for i := 0; i < 10; i++ { // Limit concurrency to prevent excessive resource usage
        wg.Add(1)
        go func() {
            defer wg.Done()
            for number := range numbers {
                if isPrime(number) {
                    mu.Lock()
                    primes = append(primes, number)
                    fmt.Printf("Found Parallel Prime: %d\n", number) // Print found primes
                    mu.Unlock()
                }
            }
        }()
    }
}
```

What do you expect?

Performance Metrics!

Start Time	End Time	Duration	Spanid	OperationName	TraceID	ThreadID	WAITBLOCKED	RUN	RUNSYSTEMCALL	WAITCPU
10:51:16.862 378 440	10:51:16.862 406 586	28,146	8409096035200089774	calculateFibonacciParallel	2881530242770011486	2410233	0	11378 19128	0	0
10:51:16.862 313 417	10:51:16.862 464 834	151,417	8717727913264033703	calculateFibonacciParallel	2881530242770011486	2410227	53358	35115 79117	9567	9567
10:51:16.862 346 033	10:51:16.862 506 809	160,776	9169217551823845228	calculateFibonacciParallel	2881530242770011486	2410227	53358	28973 70736	9567	9567
10:51:16.862 278 217	10:51:16.862 517 355	239,138	2881530242770011486	fibonacci	2881530242770011486	2410227	53358	66529 115698	9567	9567
10:51:21.713 643 539	10:51:21.713 676 923	33,384	6514851310001073122	calculateFibonacciSequenc	5476396147119877088	2410540	0	9424 25903	0	0
10:51:21.713 692 777	10:51:21.713 712 263	19,486	7465671310475914216	calculateFibonacciSequenc	5476396147119877088	2410540	0	5235 16687	0	0
10:51:21.713 632 643	10:51:21.713 729 933	97,290	8807225039319203146	calculateFibonacciSequenc	5476396147119877088	2410540	0	33155 66468	0	0
10:51:21.713 609 456	10:51:21.713 754 727	145,271	5476396147119877088	fibonacci	5476396147119877088	2410540	0	52978 94605	0	0

Trace	Timestamp	Channel	CPU	Event type	Contents
kernel	10:51:16.862 308 668	kernelchannel_5	5	sched_switch	prev_comm=swapper/5, prev_tid=0, prev_prio=20, prev_state=0, next_comm=fibonacci-par
kernel	10:51:16.862 309 576	kernelchannel_14	14	writeback_dirty_page	name=259:0, ino=21764999, index=2006, context.packet_seq_num=21, context.cpu_id=14, c
kernel	10:51:16.862 309 576	kernelchannel_8	8	power_cpu_idle	state=2, cpu_id=8, context.packet_seq_num=10, context.cpu_id=8, context_tid=0, cont
kernel	10:51:16.862 311 042	kernelchannel_14	14	writeback_mark_inode_dirty	name=259:0, ino=21764999, state=7, flags=4, context.packet_seq_num=21, context.cpu_id=
kernel	10:51:16.862 311 042	kernelchannel_4	4	syscall_exit_clock_gettime	ret=0, tp=140730733719008, context.packet_seq_num=8, context.cpu_id=4, context_tid=24
kernel	10:51:16.862 311 042	kernelchannel_5	5	syscall_exit_futex	ret=0, uaddr=824634769736, uaddr2=0, context.packet_seq_num=10, context.cpu_id=5, con
kernel	10:51:16.862 312 439	kernelchannel_1	1	syscall_entry_recvmmsg	fd=33, msg=140722553225216, flags=0, context.packet_seq_num=7, context.cpu_id=1, conb
kernel	10:51:16.862 312 439	kernelchannel_14	14	writeback_mark_inode_dirty	name=259:0, ino=21764999, state=7, flags=7, context.packet_seq_num=21, context.cpu_id=
kernel	10:51:16.862 312 439	kernelchannel_4	4	syscall_entry_clock_gettime	which_clock=1, context.packet_seq_num=8, context.cpu_id=4, context_tid=2410227, conte
kernel	10:51:16.862 312 439	kernelchannel_5	5	syscall_entry_clock_gettime	which_clock=1, context.packet_seq_num=10, context.cpu_id=5, context_tid=2410233, conte
kernel	10:51:16.862 313 417	kernelchannel_14	14	writeback_dirty_inode_start	name=259:0, ino=21764999, state=7, flags=7, context.packet_seq_num=21, context.cpu_id=
kernel	10:51:16.862 313 417	userchannel_4	4	jaeger_ust_start_span	trace_id_high=0, trace_id_low=2881530242770011486, span_id=8717727913264033703, par
kernel	10:51:16.862 314 814	kernelchannel_1	1	kmem_kfree	call_site=0xffffffff9aaa2cd, ptr=0x0, context.packet_seq_num=7, context.cpu_id=1, cont
kernel	10:51:16.862 314 814	kernelchannel_14	14	kmem_cache_alloc	call_site=0xffffffff94be865, ptr=0xffff99a400c5a80, bytes_req=56, bytes_alloc=56, gfp_fi
kernel	10:51:16.862 314 814	kernelchannel_4	4	syscall_exit_clock_gettime	ret=0, tp=140730733717984, context.packet_seq_num=8, context.cpu_id=4, context_tid=24
kernel	10:51:16.862 314 814	kernelchannel_5	5	syscall_exit_clock_gettime	ret=0, tp=139832860861152, context.packet_seq_num=10, context.cpu_id=5, context_tid=2
kernel	10:51:16.862 315 792	kernelchannel_1	1	syscall_exit_recvmmsg	ret=-11, msg=140722553225216, context.packet_seq_num=7, context.cpu_id=1, context_tid
kernel	10:51:16.862 315 792	kernelchannel_14	14	block_touch_buffer	dev=271581186, sector=87032280, size=4096, context.packet_seq_num=21, context.cpu_id=
kernel	10:51:16.862 315 792	kernelchannel_5	5	syscall_entry_clock_gettime	which_clock=1, context.packet_seq_num=10, context.cpu_id=5, context_tid=2410233, conte
kernel	10:51:16.862 317 188	kernelchannel_1	1	syscall_entry_recvmmsg	fd=33, msg=140722553225232, flags=0, context.packet_seq_num=7, context.cpu_id=1, conb
kernel	10:51:16.862 317 188	kernelchannel_14	14	kmem_cache_free	call_site=0xffffffff94bf0e6, ptr=0xffff99a400c5a80, name=mbcache, context.packet_seq
kernel	10:51:16.862 318 166	kernelchannel_1	1	kmem_kfree	call_site=0xffffffff9aaa2cd, ptr=0x0, context.packet_seq_num=7, context.cpu_id=1, cont
kernel	10:51:16.862 318 166	kernelchannel_14	14	writeback_dirty_inode	name=259:0, ino=21764999, state=7, flags=7, context.packet_seq_num=21, context.cpu_id=

63

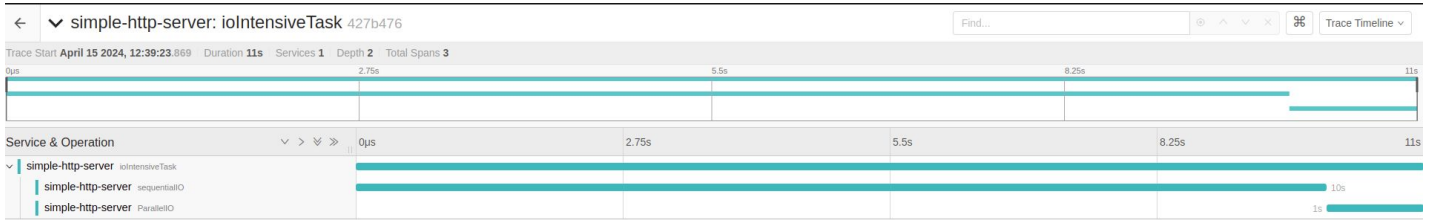
Is parallel mechanism good for all situations?

- When Parallelism Is Beneficial
 - I/O-Intensive Tasks
 - Highly Concurrent Applications
 - Tasks with High Latency Operations
- When Parallelism May Not Be Beneficial
 - Computationally Intensive Tasks
 - Memory Bandwidth Saturation
 - Scaling Issues

I/O intensive tasks!

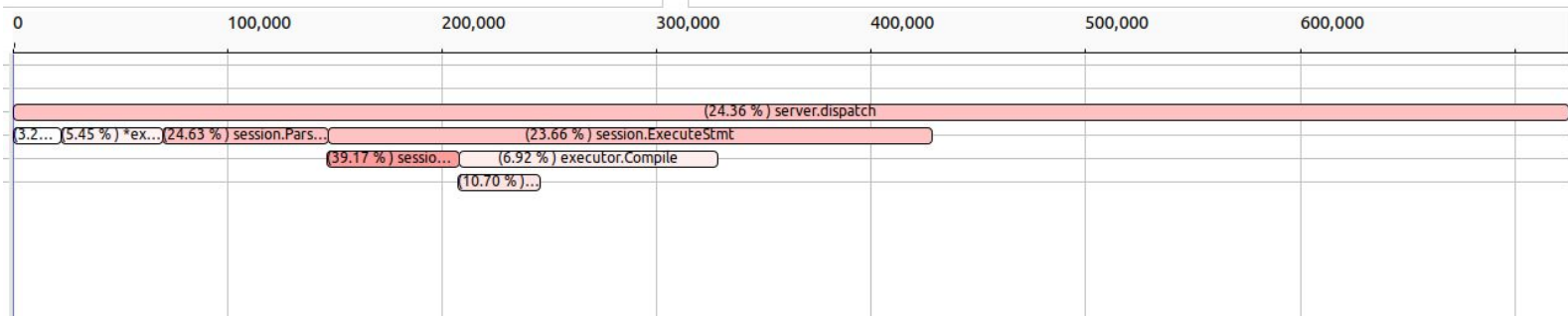
Start Time	End Time	Duration	SpanId	OperationName	TraceID	ThreadID	WAITBLOCKED	RUN	RUNSYSTEMCALL	INTERRUPTED	WAITCPU
12:40:39.048186492	12:40:49.056377540	10.008191048	7716729211719815633	sequentialIO	1543448237193233578	2606198	10006110793	566716	789417	14386	712288
12:40:49.056469731	12:40:50.057822420	1.001352689	273328434949176064	ParallelIO	1543448237193233578	2606198	1001026179	217319	99977	0	12431
12:40:39.048166936	12:40:50.057847773	11.009680837	1543448237193233578	ioIntensiveTask	1543448237193233578	2606198	11007136972	920359	910197	14386	724719

Trace	Timestamp	Channel	CPU	Event type	Contents	TID	Prio	PID
kernel	12:40:49.056355191	kernelchannel_0	0	timer_hrtimer_start	hrtimer=0xffff99aa2e423220, function=0xfffffa9188da0, expires=52591356000000, soft...	2606198	20	0
kernel	12:40:49.056356658	kernelchannel_15	15	syscall_exit_clock_gettime	ret=0, tp=140731946925616, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198	2606198	20	2606198
kernel	12:40:49.056359102	kernelchannel_0	0	power_cpu_idle	state=3, cpu_id=0, context.packet_seq_num=22, context.cpu_id=0, context_tid=0, context_...	2606198	20	0
kernel	12:40:49.056359102	kernelchannel_15	15	syscall_entry_futex	uaddr=824634769736, op=129, val=1, utime=0, uaddr2=0, val3=0, context.packet_seq_num=2606198	2606198	20	2606198
kernel	12:40:49.056362874	kernelchannel_15	15	sched_waking	comm=ioIntensivePara, tid=2606204, prio=20, target_cpu=13, context.packet_seq_num=17, 2606198	2606198	20	2606198
kernel	12:40:49.056366156	kernelchannel_15	15	syscall_exit_futex	ret=1, uaddr=824634769736, uaddr2=0, context.packet_seq_num=17, context.cpu_id=15, co...	2606198	20	2606198
kernel	12:40:49.056376074	kernelchannel_15	15	syscall_entry_clock_gettime	which_clock=1, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198, conl...	2606198	20	2606198
ust/uid/1000/64-bit	12:40:49.056377540	userchannel_15	15	jaeger_ust_end_span	trace_id_high=0, trace_id_low=1543448237193233578, span_id=7716729211719815633, dur...	2606198	20	2606198
kernel	12:40:49.056378937	kernelchannel_15	15	syscall_exit_clock_gettime	ret=0, tp=140731946924976, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198	2606198	20	2606198
kernel	12:40:49.056380668	kernelchannel_15	15	syscall_entry_clock_gettime	which_clock=0, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198, conl...	2606198	20	2606198
kernel	12:40:49.056390042	kernelchannel_15	15	syscall_exit_clock_gettime	ret=0, tp=140731946925904, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198	2606198	20	2606198
kernel	12:40:49.056391998	kernelchannel_15	15	syscall_entry_clock_gettime	which_clock=1, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198, conl...	2606198	20	2606198
kernel	12:40:49.056394791	kernelchannel_15	15	syscall_exit_clock_gettime	ret=0, tp=140731946925888, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198	2606198	20	2606198
kernel	12:40:49.056417001	kernelchannel_15	15	kmem_mm_page_alloc	page=0xffffcf9112e1f80, pfn=4503678, order=0, gfp_flags=17829322, migratetype=1, cont...	2606198	20	2606198
kernel	12:40:49.056425103	kernelchannel_15	15	kmem_mm_page_alloc	page=0xffffcf90a68a880, pfn=2728610, order=0, gfp_flags=17829322, migratetype=1, cont...	2606198	20	2606198
kernel	12:40:49.056457858	kernelchannel_15	15	kmem_mm_page_alloc	page=0xffffcf9151f1e80, pfn=5336890, order=0, gfp_flags=17829322, migratetype=1, cont...	2606198	20	2606198
kernel	12:40:49.056461630	kernelchannel_15	15	kmem_mm_page_alloc	page=0xffffcf90533ea80, pfn=1396650, order=0, gfp_flags=17829322, migratetype=1, cont...	2606198	20	2606198
kernel	12:40:49.056468335	kernelchannel_15	15	syscall_entry_clock_gettime	which_clock=1, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198, conl...	2606198	20	2606198
ust/uid/1000/64-bit	12:40:49.056469731	userchannel_15	15	jaeger_ust_start_span	trace_id_high=0, trace_id_low=1543448237193233578, span_id=273328434949176064, pare...	2606198	20	2606198
kernel	12:40:49.056471128	kernelchannel_15	15	syscall_exit_clock_gettime	ret=0, tp=140731946924864, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198	2606198	20	2606198
kernel	12:40:49.056475389	kernelchannel_15	15	syscall_entry_clock_gettime	which_clock=0, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198, conl...	2606198	20	2606198
kernel	12:40:49.056478252	kernelchannel_15	15	syscall_exit_clock_gettime	ret=0, tp=140731946925808, context.packet_seq_num=17, context.cpu_id=15, context_tid=2606198	2606198	20	2606198



What we created!

Differential Flame Graph and System metrics



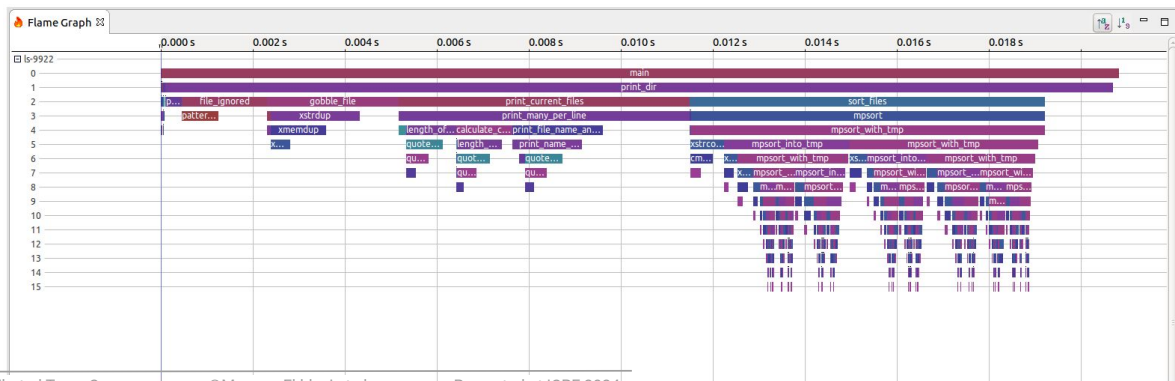
Start Time	End Time	Duration	SpanId	OperationName	TraceID	ThreadID	RUN	RUNSYSTEMCALL	INTERRUPTED	WAITCPU	WAITUNKNOWN	WAITFORK	UNKNOWN
11:06:56.933 952 078	11:06:56.933 980 155	28,077	4442450064020525663	*executor.ShowExec.Next	1618081973820362342	1805943	17108	12915	0	0	0	0	0
11:06:56.934 022 060	11:06:56.934 042 663	20,603	6023047575761396435	*executor.ShowExec.Next	1618081973820362342	1805943	7188	15709	0	0	0	0	0
11:07:46.711 230 210	11:07:46.711 249 416	19,206	2871436018761349651	*executor.ShowExec.Next	7968810417120810813	1813228	10054	11450	0	0	0	0	0
11:07:46.711 292 718	11:07:46.711 303 334	10,616	1574938190389480438	*executor.ShowExec.Next	7968810417120810813	1813228	4326	8727	0	0	0	0	0
11:06:56.933 723 277	11:06:56.933 843 056	119,779	7456353795818001242	executor.Compile	1618081973820362342	1805943	68216	53966	0	0	0	0	0
11:07:46.711 033 606	11:07:46.711 132 920	99,314	2159668538522373865	executor.Compile	7968810417120810813	1813228	52290	49358	0	0	0	0	0
11:06:56.933 567 600	11:06:56.933 592 324	24,724	3369336188777566051	Search: Duplicate	1618081973820362342	1805943	10539	16130	0	0	0	0	0
11:07:46.710 746 696	11:07:46.710 772 468	25,772	8310292673920133684	Search: Duplicate	7968810417120810813	1813228	10960	16687	0	0	0	0	0

Why differential flame graph

- Software performance analysis
- Latency detection
- New Software version or system updates
- Runtime discrepancies
- Pinpoint root causes

Flame Graph

- An aggregated view of the function calls from the call stack
- Each entry (box) represents a function in the stack.
- The x-axis represents total duration (execution time) and not absolute time.
- The width of an entry is the total time spent in that function, including time spent calling the children.
- Comprehensive yet simple and concise view, ideal for comparison



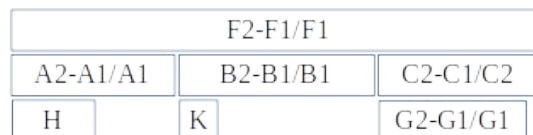
DTraComp: Distributed Trace Compare ©Maryam Ekhlesi et al. Presented at ICPE 2024

Differential Flame Graph-Computation

Suppose having two traces, we need to compare them



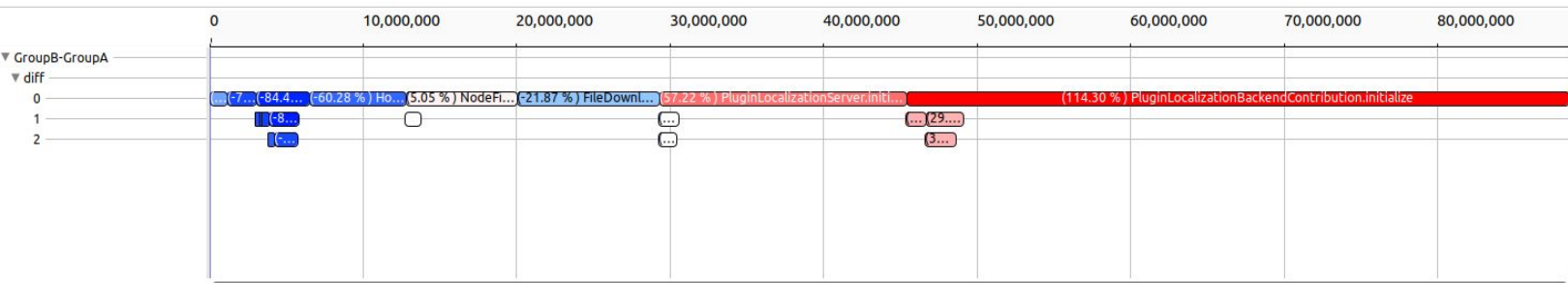
Differential Flame Graph



Differential Flame Graph-Representation

- Color

- Flame Graph: The color of each function is randomly assigned based on the function name
- Differential Flame Graph: colors are assigned based on the difference ratio
 - Red color palette for longer execution
 - Blue color palette for shorter execution
 - White for difference ratio between -5% and 5%

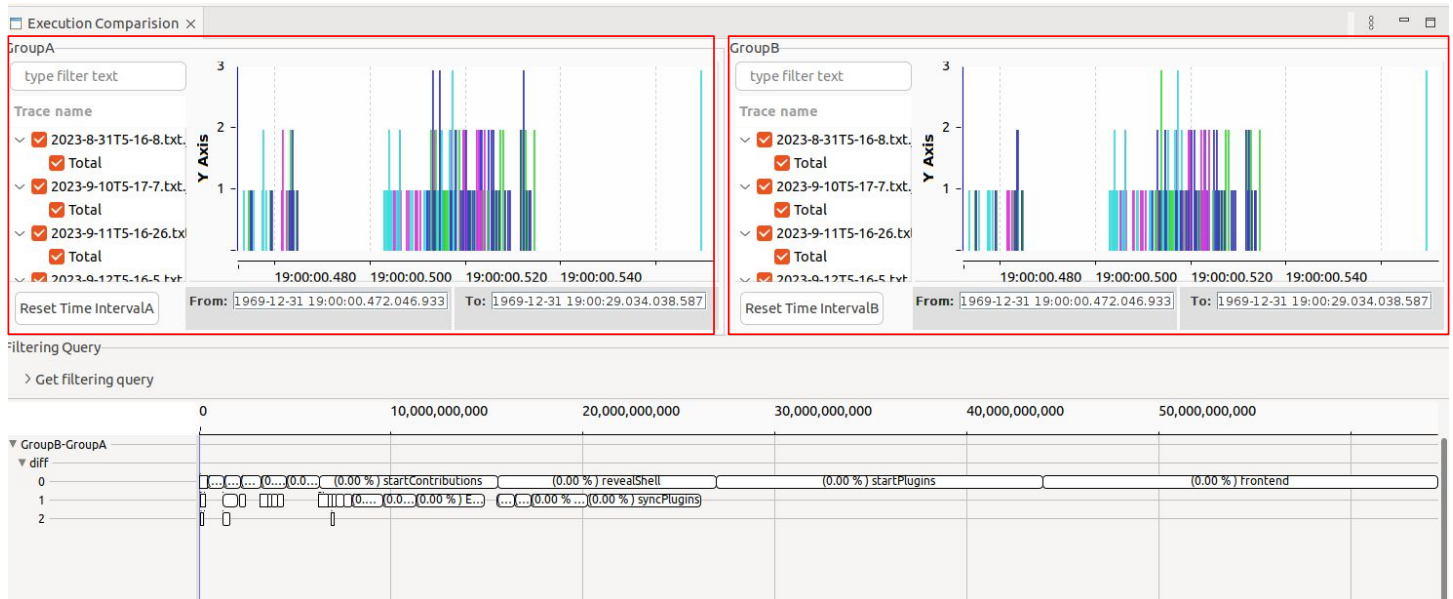


Execution Comparison

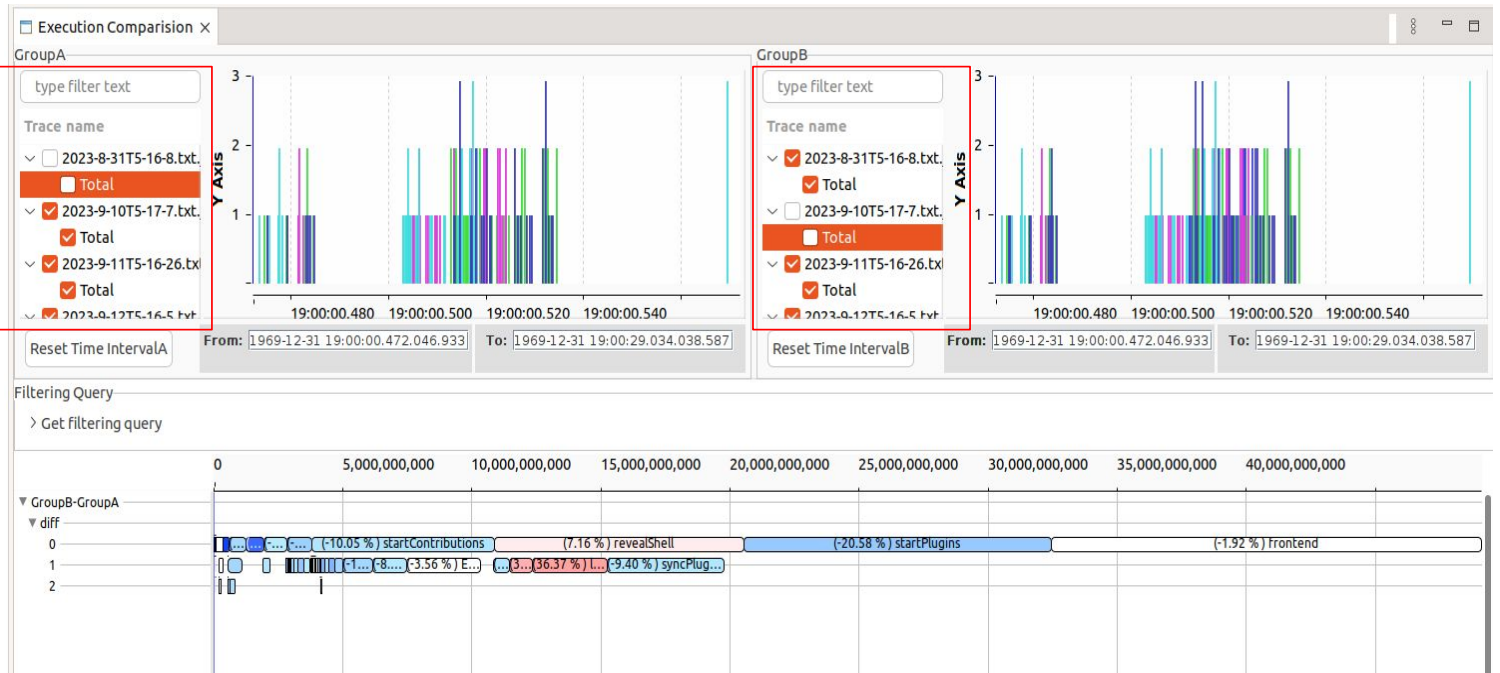
- We have two mode for execution comparison
 - Sequential traces
 - Distributed microservice-based system's traces (request comparison)

Sequential trace comparison

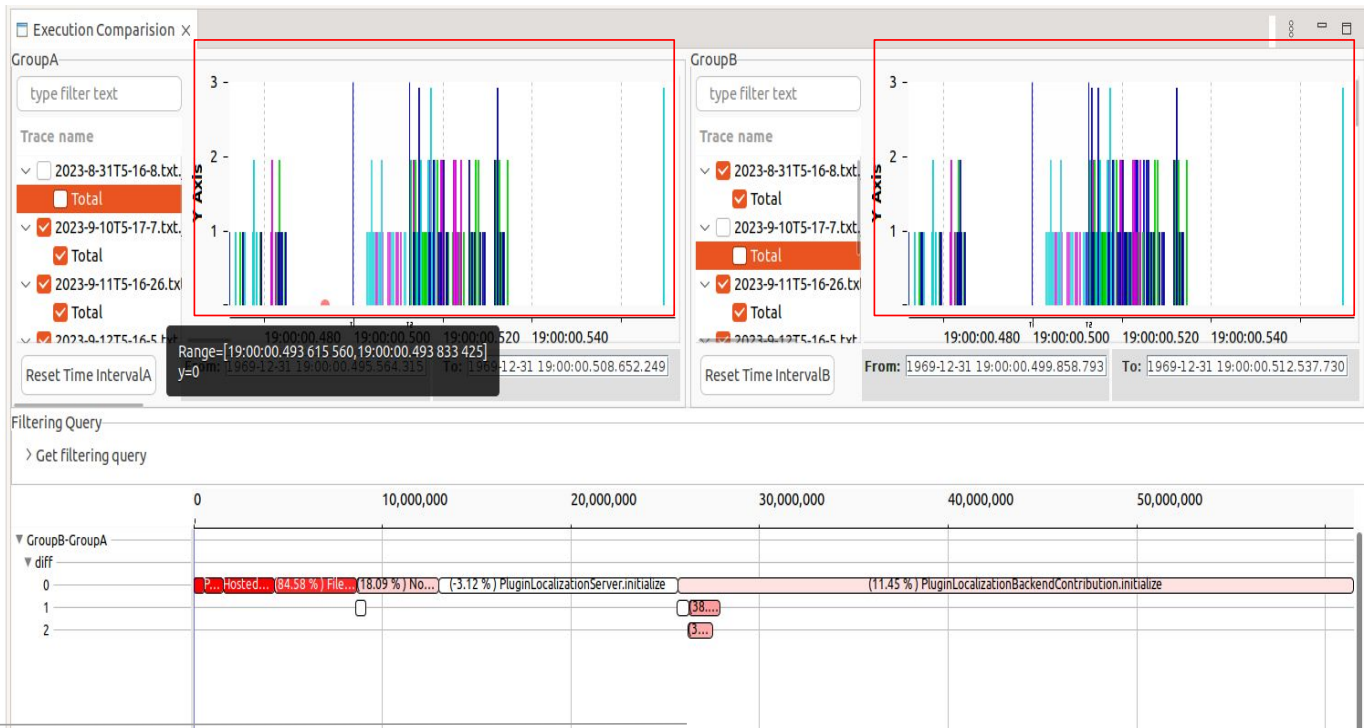
Select two groups to be compared



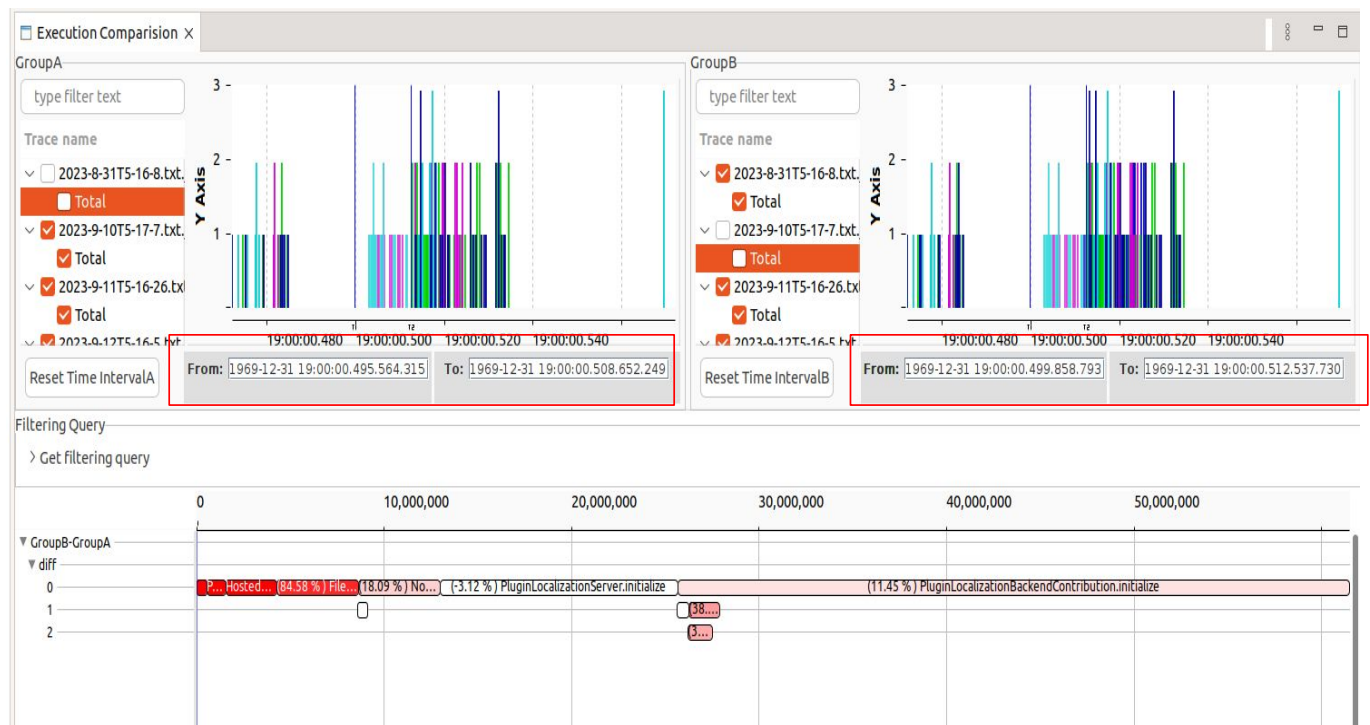
Select Relevant Traces



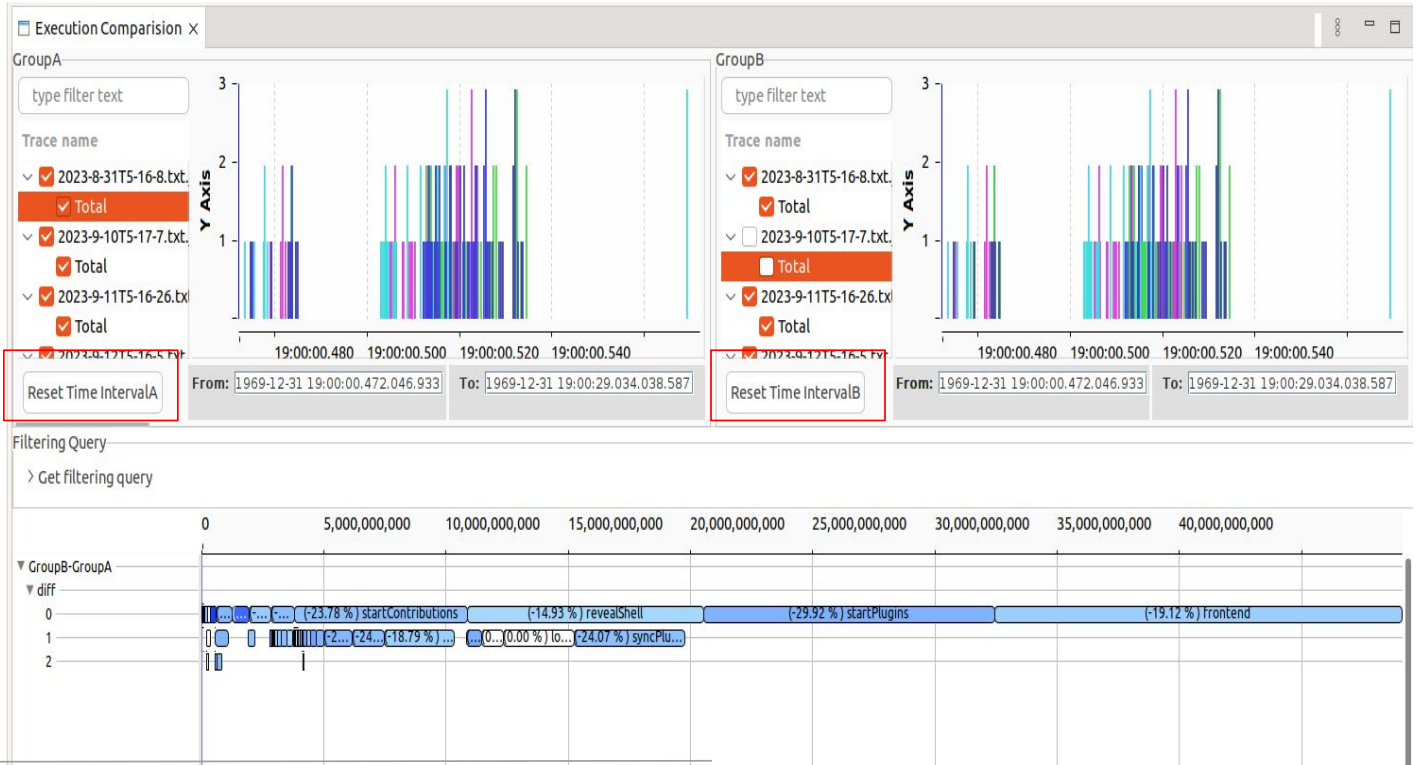
Select desired time range graphically



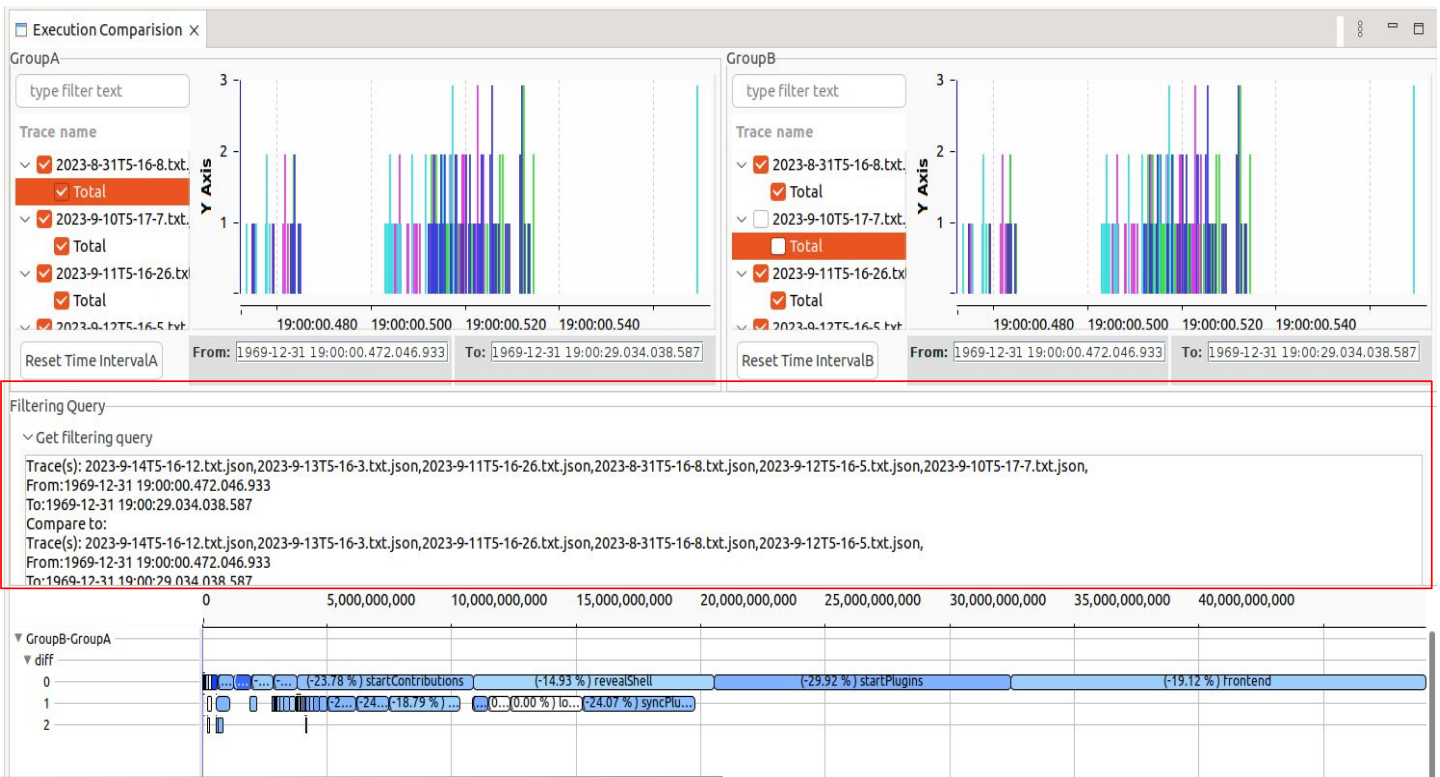
Select desired time range textually



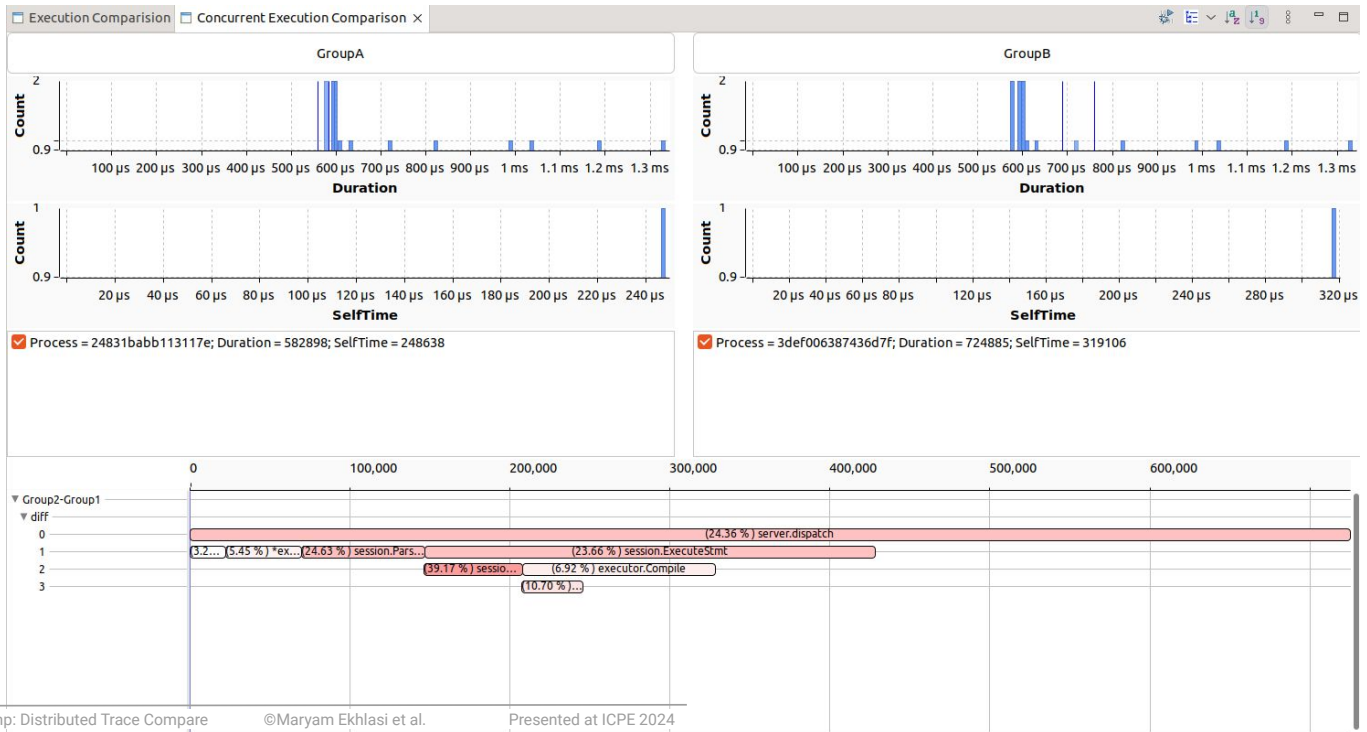
Reset to initial state



Experiment reproducibility



Request comparison



Let's collect traces and see the result in TraceCompass!

Find all the commands here

<https://github.com/maryamekhlasi/ICPE2024>



Installing LTTng on Ubuntu

1. Add the LTTng Stable 2.13 PPA repository and update the list of packages

```
# apt-add-repository ppa:lttng/stable-2.13  
# apt-get update
```

2. Install the main LTTng 2.13 packages

```
# apt-get install lttng-tools  
# apt-get install lttng-modules-dkms  
# apt-get install liblttng-ust-dev
```

Recording Kernel Events With LTTng

1. Checking the status of session daemon

```
$ systemctl status lttng-sessiond
```

2. Starting the session

```
$ sudo systemctl start lttng-sessiond
```

3. Create tracing session

```
$ sudo lttng create ICPE2024 --output="/home/maryam/Poly/Dorsal/traces/dataset/ICPE202"
```

4. Enable Events

```
$ sudo lttng enable-event -k -a
```

5. Start Tracing

```
$ sudo lttng start
```

6. Do something! Open a browser!

```
$ wget www.google.com
```

7. Stop and destroy the Tracing session

```
$ sudo lttng destroy
```

83

Installing TraceCompass

1. Installing Java

```
java -version  
sudo apt update  
sudo apt install default-jre
```

2. Installing TraceCompass

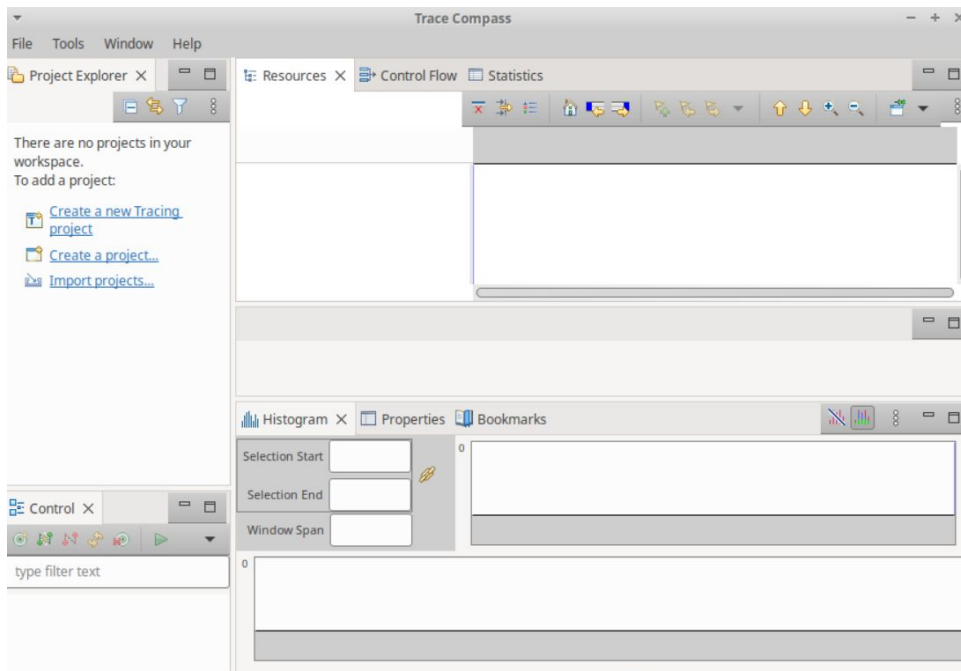
Go to <https://eclipse.dev/tracecompass/>

```
Tar xf <name>
```

```
Cd trace-compass
```

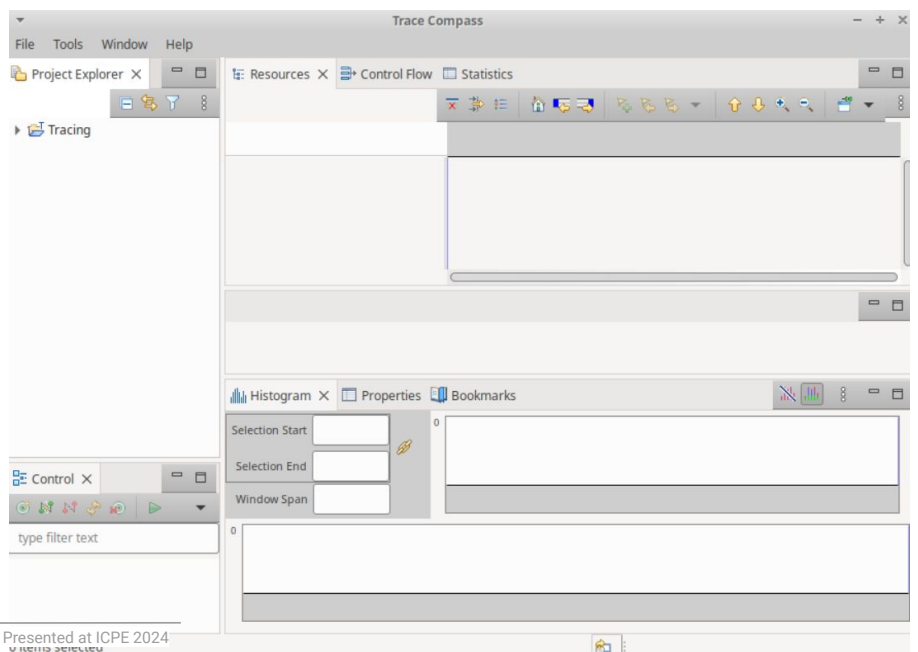
```
./tracecompass
```

Empty space TraceCompass

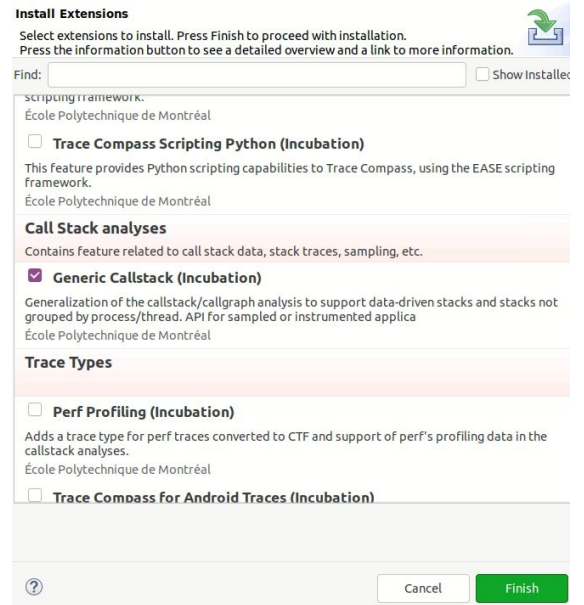
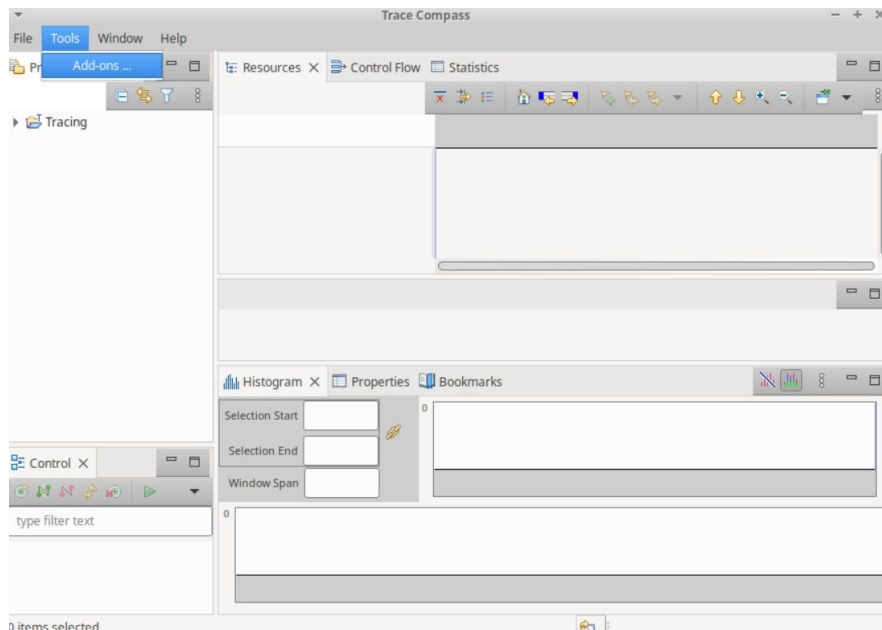


Create a tracing project

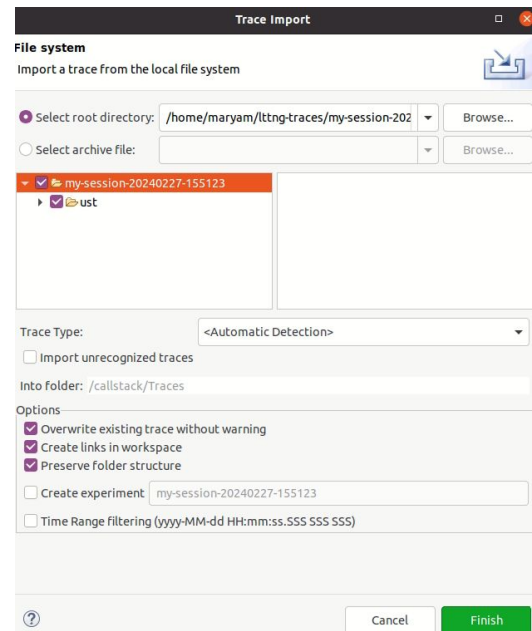
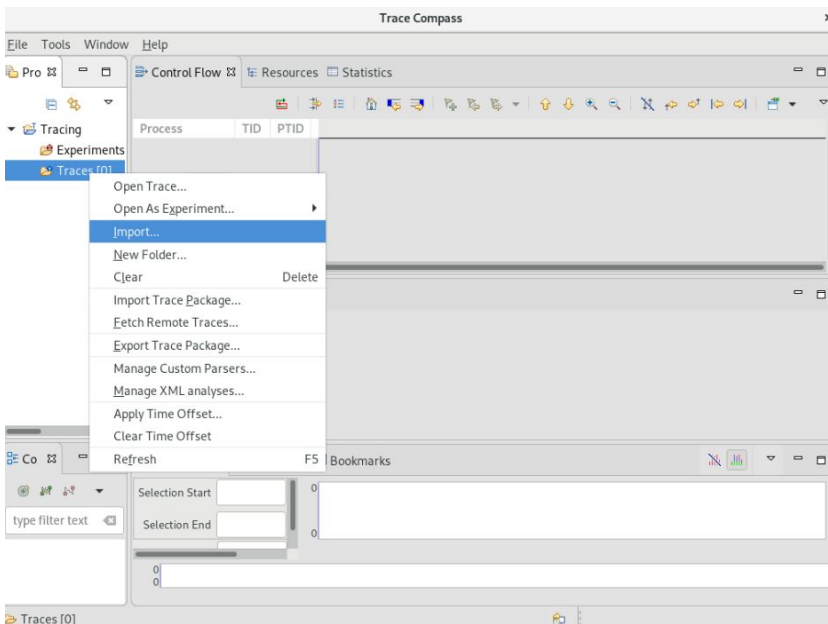
File >> new tracing project



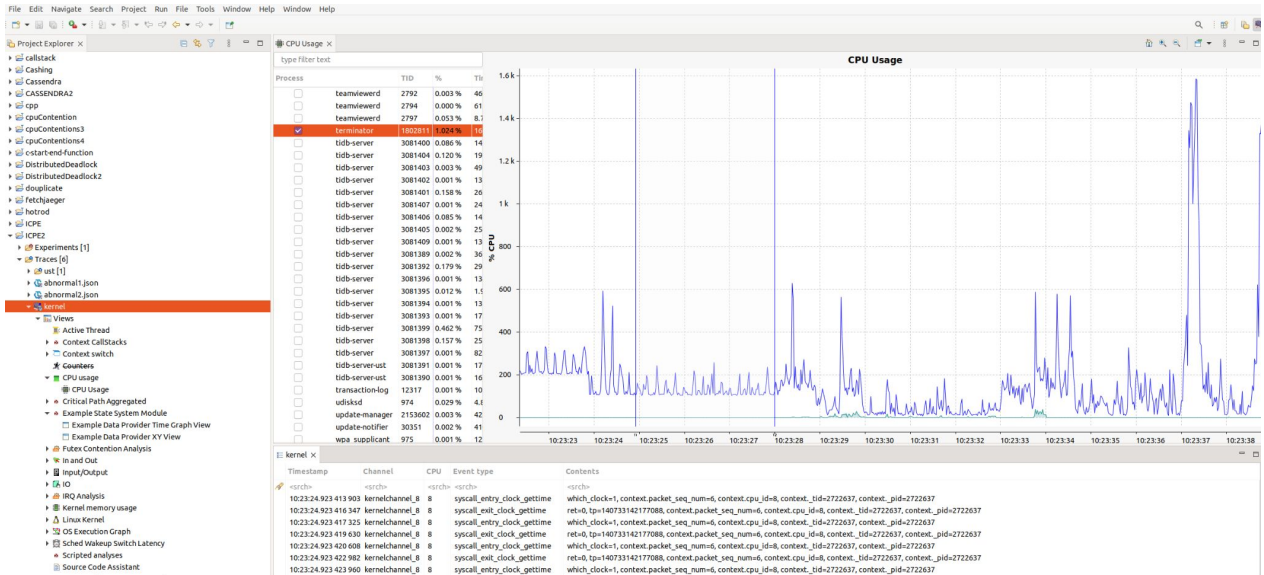
Install the required add-ons



Import traces



View results

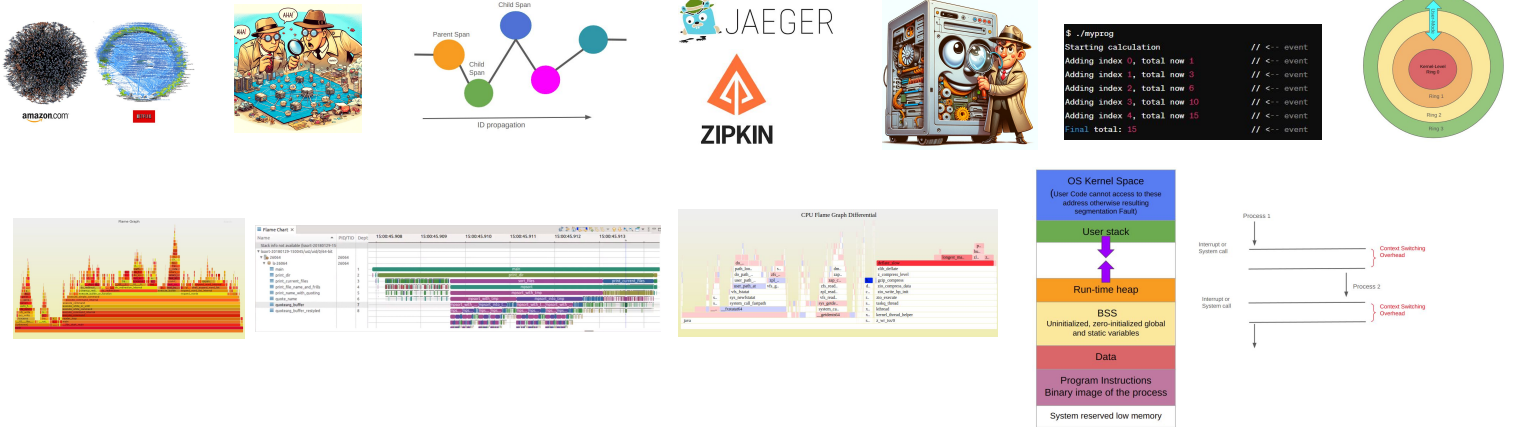


Source codes are available for free!



1. <https://github.com/eclipse-tracecompass-incubator/org.eclipse.tracecompass.incubator>
2. <https://git.eclipse.org/r/c/tracecompass.incubator/org.eclipse.tracecompass.incubator/+196496>

Thank you!



References

1. <https://lttn.org/docs/v2.13/>
2. <https://eclipse.dev/tracecompass/>
3. <https://github.com/dorsal-lab/Tracevizlab/>
4. <https://www.brendangregg.com/>
5. Maryam Ekhlesi, Fatemeh Faraji Daneshgar, Michel Dagenais, et al. DTraComp: Comparing distributed execution traces for understanding intermittent latency sources. Authorea. October 18, 2023.