

ไลฟ์ไมเกรชันแบบเทอร์คบนเวอร์ชวลแมชชีนด้วยเทคนิคการเชื่อมต่อเครือข่ายแบบขนาน

Thread-based Live Migration of Virtual Machines Using Parallel Network Connections

พิทักษ์ แทนแก้ว¹ และ กษิดิศ ชาญเชียว²

ภาควิชาวิทยาการคอมพิวเตอร์ คณะวิทยาศาสตร์และเทคโนโลยี มหาวิทยาลัยธรรมศาสตร์

99 หมู่ 18 ถ.พหลโยธิน ต.คลองหนึ่ง อ.คลองหลวง จ.ปทุมธานี 12121 ประเทศไทย

Email: ¹5209035194@student.cs.tu.ac.th, ²kasidit@cs.tu.ac.th

บทคัดย่อ

บทความนี้นำเสนอวิธีการไลฟ์ไมเกรชันของเวอร์ชวลแมชชีนแบบใหม่เรียกว่า ทีเอลเอ็ม (TLM) ที่ทำให้การเคลื่อนย้ายเวอร์ชวลแมชชีนระหว่างเครื่องคอมพิวเตอร์จริงสองเครื่องเป็นไปได้อย่างมีประสิทธิภาพ วิธีการทีเอลเอ็มนำเสนอการออกแบบใหม่สองประการได้แก่ 1) การใช้เทอร์คเพื่อส่งสถานะของเวอร์ชวลแมชชีนไปยังเครื่องคอมพิวเตอร์ปลายทางในขณะที่เวอร์ชวลแมชชีนประมวลผล และ 2) มีการส่งข้อมูลแบบขนานเพื่อให้การเคลื่อนย้ายเวอร์ชวลแมชชีนเร็วขึ้น ระบบต้นแบบของทีเอลเอ็มได้รับพัฒนาขึ้นบนโอเพ่นซอร์สซอฟต์แวร์เคอร์เนล (Kernel-based Virtual Machine) และทดสอบกับเวอร์ชวลแมชชีนขนาดใหญ่ที่กำลังประมวลผล NAS Parallel Benchmark ผู้วิจัยพบว่าทีเอลเอ็มมีประสิทธิภาพมากกว่าวิธีการไมเกรชันเดิมของเคอร์เนลทั้งแบบไลฟ์ไมเกรชันและแบบหยุดแล้วไมเกรชันหลายเท่า

คำสำคัญ: เวอร์ชวลแมชชีน ไลฟ์ไมเกรชัน การเชื่อมต่อเครือข่ายแบบขนาน

Abstract

This paper describes a novel Thread-based Live Migration (TLM) of virtual machines where additional threads are created to facilitate the state transfer of a virtual machine over parallel TCP connections concurrently with the virtual machine's execution. We have implemented TLM on the Open Source Kernel-based Virtual Machine (KVM) software and conducted preliminary experiments by migrating a large SMP virtual machine running computation and memory intensive NAS Parallel Benchmarks. The results indicate significantly better performance of TLM over the traditional live and non-live migration mechanisms.

keywords: virtual machine, live migration, parallel network connections

1. บทนำ

ในระบบประมวลผลแบบกลุ่มเมฆ (Cloud Computing) ซึ่งประกอบไปด้วย การใช้งานเวอร์ชวลแมชชีนจำนวนมาก ความสามารถในการไลฟ์ไมเกรชัน หรือการเคลื่อนย้ายเวอร์ชวลแมชชีนในขณะที่เวอร์ชวลแมชชีนกำลังประมวลผลและมีการจัดจ้งหะการทำงานของเวอร์ชวลแมชชีนน้อยที่สุดนั้นเป็นความสามารถที่เป็นประโยชน์หลายประการ อาทิ

- เพิ่มประสิทธิภาพของแอปพลิเคชันที่ปฏิบัติงานอยู่บนเวอร์ชวลแมชชีนโดยไมเกรชันไปยังเครื่องคอมพิวเตอร์จริงที่มีประสิทธิภาพสูงขึ้น หรือไมเกรชันไปยังเครื่องที่อยู่ใกล้กับฐานข้อมูล
 - การเพิ่มความทนทานต่อความผิดพลาด (Fault Tolerance) โดยไมเกรชันจากเครื่องคอมพิวเตอร์ที่มีอุปกรณ์บางส่วนเสียหายไปยังเครื่องที่ทำงานปกติ
 - ช่วยการบริหารการใช้งานทรัพยากรภายในคลัสเตอร์ให้มีประสิทธิภาพโดยไมเกรชันเวอร์ชวลแมชชีนที่อาจกระจายอยู่บนเครื่องคอมพิวเตอร์หลาย ๆ เครื่องให้ไปอยู่บนเครื่องเดียวเพื่อประหยัดพลังงาน
- อย่างไรก็ตาม วิธีการไลฟ์ไมเกรชันในปัจจุบันยังมีข้อจำกัดและไม่สามารถปฏิบัติงานได้อย่างมีประสิทธิภาพเมื่อเวอร์ชวลแมชชีนเป็นแบบ SMP ที่ใช้ซีพียูคอร์แบบขนานและมีขนาดของหน่วยความจำขนาดใหญ่ซึ่งเป็นลักษณะของเวอร์ชวลแมชชีนที่มีให้บริการกับลูกค้าตามปกติในคลัสเตอร์ [1] ดังที่แสดงในตารางที่ 1 ซึ่งแบบ Extra Large และ Double และ Quadruple Extra Large นั้นใช้สำหรับแอปพลิเคชันที่เน้นการใช้งานฐานข้อมูลและการแคชข้อมูลในหน่วยความจำ ส่วนแบบ Cluster นั้นใช้ในการประมวลผลสมรรถนะสูง
- การทำไลฟ์ไมเกรชันแบบเดิมจะใช้เวลาอย่างมากหรืออาจใช้ไม่ได้เลยสำหรับเวอร์ชวลแมชชีนที่มีปริมาณการประมวลผลมากและมีขนาดของหน่วยความจำมาก ดังนั้นการไมเกรชันจึงต้องหันไปใช้วิธีหยุดการทำงานของเวอร์ชวลแมชชีนแล้วส่งสถานะไปยังคอมพิวเตอร์ปลายทางซึ่งเป็นทางเลือกที่ยอมรับได้สำหรับเวอร์ชวลแมชชีนที่รัน Batch Jobs แต่ไม่เหมาะสมกับเวอร์ชวลแมชชีนที่ทำงานเป็นฐานข้อมูลหรือแคชข้อมูลหรือแอปพลิเคชันที่ต้องการการตอบสนอง

ตารางที่ 1 ชนิดของเวอร์ชวลแมชชีนของบริษัท Amazon

Instant Types	No. of Virtual	Memory Size	Disk Space
	Core	(Gigabytes)	(Gigabytes)
Large	2	7.5	850
Extra Large	4	15	1,690
Double Extra Large	4	34.2	850
Quadruple Extra Large	8	68.4	1,690
Cluster Quad Extra Large	8	23	1,690
Cluster Eight Extra Large	16	60.5	3,370

งานวิจัยนี้นำเสนอการทำไลฟไทม์ไมเกรชันแบบใหม่ที่สามารถเคลื่อนย้ายเวอร์ชวลแมชชีนแบบ SMP ขนาดใหญ่ได้ในขณะที่ทำงานและมีการจัดจังหวะการทำงานของเวอร์ชวลแมชชีนน้อยกว่าการไมเกรชันแบบเดิมหลายเท่า โดยนำที่เอลซี (TLC) [7] ซึ่งใช้เวอร์ชวลไลเซชันซอฟต์แวร์ที่เป็นโอเพ่นซอร์สชื่อควีเอ็ม มาพัฒนาต่อ ที่เอลซีนั้นใช้วิธีสร้างเทรคขึ้นมาเพิ่มเพื่อทำหน้าที่เช็คพอยน์ โดยเวอร์ชวลแมชชีนและเทรคดังกล่าว จะทำงานไปพร้อม ๆ กัน

สำหรับวิธีการใหม่หรือที่เอลเอ็ม ได้รับการออกแบบให้มีความยืดหยุ่น เทรคแรกทำหน้าที่สแกนหน่วยความจำและส่งไปยังโฮสปลายทาง เทรคที่สองเลือกส่งเฉพาะข้อมูลที่เกิดขึ้นใหม่ เทรคที่สามคอยส่งข้อมูลที่เหลือ โดยเทรคแรกและเทรคที่สองจะส่งข้อมูลในลักษณะขนานกันไป

จากการทดลองกับเวอร์ชวลแมชชีนที่ประมวลผลโปรแกรม mg Class D และ sp Class D ของ NAS Parallel Benchmarks พบว่าที่เอลเอ็มสามารถทำไลฟไทม์ไมเกรชันได้ภายในเวลาจำกัดและสามารถวัดความคืบหน้าได้ ในขณะที่วิธีการเดิมใช้เวลานานมากกว่า Benchmark จะทำงานจบจึงไมเกรท นอกจากนี้ที่เอลเอ็มมีระยะเวลาการจัดจังหวะน้อยกว่าการไมเกรชันแบบหยุดแล้วจึงไมเกรทเดิมของควีเอ็มถึง 7 เท่าและ 5.2 เท่าสำหรับโปรแกรม mg Class D และ sp Class D

จากผลการทดลองเบื้องต้นพบว่า ที่เอลเอ็ม ซึ่งมีการส่งข้อมูลแบบขนานในลักษณะพื้นฐาน ให้ผลที่ดีกว่างานอื่น ๆ ก่อนหน้านี้ ผู้วิจัยมีความคิดว่า การเพิ่มการเชื่อมต่อในการส่งข้อมูลให้มากขึ้น น่าจะทำให้การทำงานของไลฟไทม์ไมเกรชันเร็วขึ้นอีก

ผู้วิจัยจึงทำการพัฒนาไลบรารีสำหรับการเชื่อมต่อเครือข่ายแบบขนานขึ้นมา เพื่อทดสอบสมมติฐานข้างต้น และได้ออกแบบการทดลองต่างหาก ด้วยการส่งข้อมูลโดยใช้การเชื่อมต่อแบบต่าง ๆ เพื่อทดสอบความเป็นไปได้ ศึกษาแนวโน้ม ก่อนจะนำไปประยุกต์ใช้กับที่เอลเอ็ม ให้มีการทำงานแบบขนานมากขึ้น ซึ่งเป็นงานที่จะทำต่อไปในอนาคต

2. ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

การทำไลฟไทม์ไมเกรชันของเวอร์ชวลแมชชีนนั้นได้รับการพัฒนาขึ้นบนไฮเปอร์ไวเซอร์ได้แก่ Xen [2] และ Vmotion ของ VMware [3] และ KVM [4] ซึ่งมีหลักการที่คล้ายกันคือการไมเกรชันแบบ Pre-Copy ซึ่งจะ

ทยอยส่งข้อมูลเป็นระยะจากเวอร์ชวลแมชชีนต้นทางไปยังปลายทางโดยที่ในการส่งเพียงหน่วยความจำแต่ละครั้งจะส่งไม่มากเพื่อไม่ให้ขัดจังหวะการทำงานของเวอร์ชวลแมชชีนมากเกินไปและทำเช่นนี้ต่อไปเรื่อย ๆ จนกระทั่งจำนวน Dirty เพจที่จำเป็นต้องส่งเหลืออยู่น้อยแล้วจึงหยุดเวอร์ชวลแมชชีนเพื่อส่งสถานะที่เหลือไปทั้งหมดแล้วโอนการทำงานไปยังเวอร์ชวลแมชชีนปลายทาง ปัญหาของวิธีการนี้คือถ้าเวอร์ชวลแมชชีนเป็นแบบ SMP และมีการเขียนหน่วยความจำเป็นปริมาณมากวิธีนี้จะต้องใช้เวลานานมากจนกว่าจะเหลือ Dirty เพจน้อยแล้วจบได้

งานวิจัยทางด้านไมเกรชันมีจุดมุ่งหมายเดียวกันคือการพยายามพัฒนาวิธีการที่ทำให้การไมเกรชันมีประสิทธิภาพและมีการจัดจังหวะการทำงานของเวอร์ชวลแมชชีนที่น้อยที่สุดเท่าที่จะทำได้ งานวิจัยของ [5] ใช้เทคนิค Delta Compression เพื่อลดขนาดของข้อมูลและเพิ่มประสิทธิภาพการทำไลฟไทม์ไมเกรชันสำหรับแอปพลิเคชันสำหรับองค์กรซึ่งยังมีการใช้งานหน่วยความจำและปริมาณการประมวลผลน้อยกว่า NPB Benchmark ที่ใช้ในการทดลองนี้ แต่อย่างไรก็ตามการทำ Delta Compression อาจนำมาใช้เพื่อเสริมประสิทธิภาพของที่เอลเอ็มได้ในอนาคต งานวิจัย [6] เสนอวิธีการไมเกรชันแบบ Post-Copy โดยย้ายการทำงานของเวอร์ชวลแมชชีนไปยังคอมพิวเตอร์ปลายทางก่อนแล้วค่อยทยอยเคลื่อนย้ายเพียงหน่วยความจำที่จำเป็นต้องใช้ในการประมวลผล On-Demand ในกรณีที่มีความจำเป็นต้องใช้หน่วยความจำมากอย่างต่อเนื่องบนเวอร์ชวลแมชชีนแบบ SMP วิธีนี้อาจทำให้เวอร์ชวลแมชชีนต้องหยุดการทำงานเป็นเวลานานได้

3. การดำเนินการวิจัย

ประกอบด้วย การออกแบบไลฟไทม์ไมเกรชันแบบเทรค การพัฒนาไลบรารีการเชื่อมต่อแบบขนาน และงานที่จะทำต่อไปในอนาคตคือการนำสองอย่างข้างต้นมาประยุกต์ใช้งานร่วมกัน

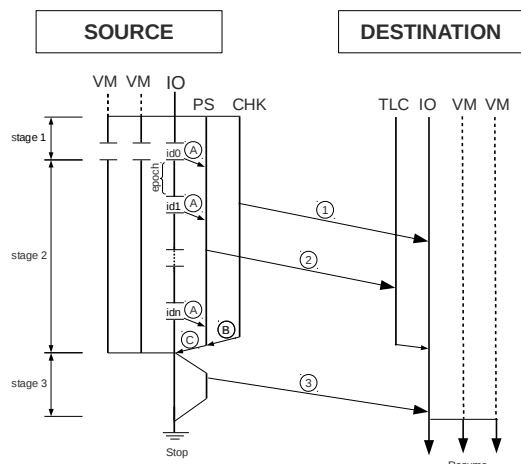
3.1 ไลฟไทม์ไมเกรชันแบบเทรค

สถาปัตยกรรมไลฟไทม์ไมเกรชันแบบเทรค (TLM - Thread-based Live Migration) ได้รับการออกแบบโดยอาศัยหลัก Pre-Copy ประกอบไปด้วยเครื่องโฮสต์ต้นทางและปลายทาง โดยที่เครื่องโฮสต์ต้นทางจะเป็นเครื่องที่รันเวอร์ชวลแมชชีน ที่จะทำไลฟไทม์ไมเกรชันไปยังเครื่องโฮสต์ปลายทางฝั่งต้นทางประกอบด้วย เทรคการทำงานของเวอร์ชวลแมชชีน (VM) เทรคไอโอ (IO) เทรคเพจเซฟเวอร์ (PS) และ เทรคเช็คพอยน์ (CHK) ฝั่งปลายทางประกอบด้วยเทรคการทำงานของเวอร์ชวลแมชชีน (VM) เทรคที่รองรับข้อมูล (IO) และเทรคที่เอลซี (TLC)

จากรูปที่ 1 เมื่อเครื่องโฮสต์ต้นทางได้รับคำสั่งให้ทำไมเกรชัน (stage 1) จะมีหยุดการทำงานของเวอร์ชวลแมชชีน ทำการกำหนดค่าเริ่มต้นต่าง ๆ สร้างเทรคเพจเซฟเวอร์และเทรคเช็คพอยน์ เทรคไอโอจะทำการตรวจสอบ dirty ที่เกิดขึ้น และแจ้งผลให้กับเทรคเพจเซฟเวอร์ (ลูกศร A)

ซึ่งใน stage 1 นี้จะยังไม่มี dirty

เมื่อเข้าสู่ stage 2 เทรคเวอร์ชวลแมชชีนและเทรคไอโอจะเริ่มทำงานต่อ เทรคเช็คพอยน์จะทำการอ่านหน่วยความจำของเวอร์ชวลแมชชีนและทยอยส่งไปยังฮอสปลายทาง (ลูกศร 1) สำหรับเทรคไอโอจะมีการหยุดทำงานเป็นช่วง ๆ ตามระยะเวลา epoch เพื่อตรวจสอบ dirty ที่เกิดขึ้นและแจ้งผลให้กับเทรคเพจเซฟเวอร์ (ลูกศร A) เพื่อให้เทรคเพจเซฟเวอร์ส่งข้อมูลเฉพาะ dirty ไปยังฮอสปลายทาง (ลูกศร 2) ทางฝั่งฮอสปลายทางมีสองเทรคที่คอยรับข้อมูล และเขียนข้อมูลลงหน่วยความจำ โดยให้ลำดับความสำคัญกับข้อมูลที่ได้รับจากเทรคเพจเซฟเวอร์มากกว่าเทรคเช็คพอยน์สำหรับกรณีที่มีการเขียนข้อมูลลงตำแหน่งเดียวกัน



รูปที่ 1 สถาปัตยกรรมไลฟไทม์ไมเกรชันแบบเทรค

เมื่อเทรคเช็คพอยน์ทำงานเสร็จจะส่งสัญญาณให้กับเทรคเพจเซฟเวอร์เพื่อให้หยุดทำงาน จากนั้นเทรคเพจเซฟเวอร์ก็จะส่งสัญญาณให้กับเทรคไอโอ เพื่อให้หยุดการทำงานของเวอร์ชวลแมชชีน เมื่อเข้าสู่ stage 3 ฮอสปลายทางจะทำการส่งข้อมูลที่เหลือไปยังฮอสปลายทาง (ลูกศร 3) เมื่อส่งครบแล้วจึงหยุดการทำงาน ทางฝั่งฮอสปลายทางเมื่อได้รับข้อมูลครบแล้วจะทำการสั่งให้เวอร์ชวลแมชชีนทำงานต่อ

3.2 ไลบรารีสำหรับการเชื่อมต่อแบบขนาน

ผู้ทำการวิจัย ได้พัฒนาไลบรารีสำหรับอำนวยความสะดวกในการเชื่อมต่อเครือข่ายแบบขนานโดยใช้โปรโตคอลที่ซีพี ตามโมเดลไคลเอนต์-เซิร์ฟเวอร์ (Client-Server Model) มีเป้าหมายเพื่อนำไปใช้กับการทำไลฟไทม์ไมเกรชันแบบเทรค จำนวนเป็น 2 ส่วนดังนี้

```
#define ADDRLEN 16
#define PORTLEN 16

typedef struct {
    char addr[ADDRLEN];
    char port[PORTLEN];
} mr_addr_port;
```

รูปที่ 2 โครงสร้างข้อมูลสำหรับการเชื่อมต่อแบบขนาน

3.2.1 โครงสร้างข้อมูล ชื่อ mr_addr_port ใช้เก็บข้อมูลชุดของหมายเลขไอพีและหมายเลขพอร์ต ของเครื่องเซิร์ฟเวอร์ โดยการใช้งานจะประกาศตัวแปรเป็นแบบอरेย์

```
mr_addr_port *mr_get_opt(const char *filename, int *n);
int *mr_get_fd_list(int n);
void mr_server_listen(int *fd, mr_addr_port *ap, int n);
int *mr_server_accept(int *fd, int n);
void mr_client_connect(int *fd, mr_addr_port *ap, int n);
```

รูปที่ 3 ฟังก์ชันสำหรับการเชื่อมต่อแบบขนาน

3.2.2 ฟังก์ชันการทำงาน ประกอบด้วยฟังก์ชันหลัก 5 ฟังก์ชัน ดังนี้

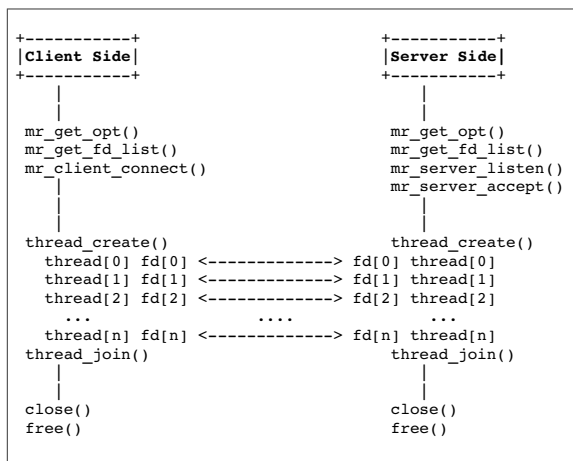
mr_get_opt(filename, n) ทำหน้าที่อ่านค่าจากไฟล์คอนฟิก ที่เก็บข้อมูลชุดของหมายเลขไอพี และหมายเลขพอร์ตของเครื่องเซิร์ฟเวอร์ เมื่อจบการทำงานจะคืนค่าเป็นอเรย์ของ mr_addr_port และ n จะถูกกำหนดให้มีค่าเท่ากับจำนวนหมายเลขไอพีและหมายเลขพอร์ต

mr_get_fd_list(n) ทำหน้าที่สร้างซ็อกเก็ตสำหรับการเชื่อมต่อ โดยคืนค่าเป็นอเรย์ของไฟล์เดสคริปเตอร์ (File Descriptor) อเรย์ที่สร้างมีขนาดเท่ากับ n

mr_server_listen(fd, ap, n) ทำการไบนด์ (Bind) หมายเลขไอพี และหมายเลขพอร์ต เข้ากับไฟล์เดสคริปเตอร์ที่ได้จาก mr_get_fd_list() และปรับสถานะของซ็อกเก็ตให้เป็นการรอรับการเชื่อมต่อ ฟังก์ชันนี้ใช้เฉพาะทางฝั่งเซิร์ฟเวอร์เท่านั้น

mr_server_accept(fd, n) รอรับการเชื่อมต่อจากไคลเอนต์ จนครบ จำนวน n การเชื่อมต่อ เมื่อจบการทำงานแล้วจะคืนค่าเป็นอเรย์ของไฟล์เดสคริปเตอร์ที่ใช้ติดต่อสื่อสารกับไคลเอนต์ ฟังก์ชันนี้ใช้เฉพาะทางฝั่งเซิร์ฟเวอร์เท่านั้น

mr_client_connect(fd, ap, n) ทำการเชื่อมต่อไปยังเซิร์ฟเวอร์จนครบจำนวน n การเชื่อมต่อ หลังจากจบการทำงานแล้ว ไคลเอนต์สามารถติดต่อสื่อสารกับเซิร์ฟเวอร์ได้ ผ่านอเรย์ของไฟล์เดสคริปเตอร์ ฟังก์ชันนี้ใช้เฉพาะทางฝั่งไคลเอนต์เท่านั้น



รูปที่ 4 รูปแบบการใช้งานไลบรารี

รูปแบบการใช้นาโบลารี่มาใช้งานจะเป็นในลักษณะมัลติเทรค ดังแสดงในรูปที่ 4

3.3 การปรับปรุงไลฟ์ไทม์เกรซันแบบเทรค

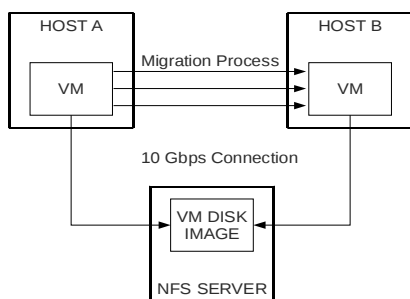
ในอนาคตผู้วิจัย จะดำเนินการปรับปรุงไลฟ์ไทม์เกรซันแบบเทรคให้มีประสิทธิภาพมากขึ้นอีก โดยมีเป้าหมายเพื่อลดระยะเวลาที่เวอร์ชวลแมชชีนหยุดการทำงาน (Downtime = Time_{stage 3}) และลดระยะเวลาทั้งหมดในกระบวนการไมเกรซัน (Elapsed Time = Time_{stage 1} + Time_{stage 2} + Time_{stage 3}) ด้วยการเพิ่มการทำงานแบบขนานให้มากขึ้น โดยนำไลบรารีที่ได้รับการพัฒนา มาประยุกต์ใช้กับไลฟ์ไทม์เกรซันแบบเทรค

4. การทดลองเบื้องต้น

การทดลองประกอบไปด้วย 2 ส่วน คือ การทดลองไลฟ์ไทม์เกรซันแบบเทรค (TLM) และ การทดลองเพื่อศึกษาความเป็นไปได้ในการใช้ไลบรารีที่พัฒนาขึ้นมา

4.1 การทดลองไลฟ์ไทม์เกรซันแบบเทรค

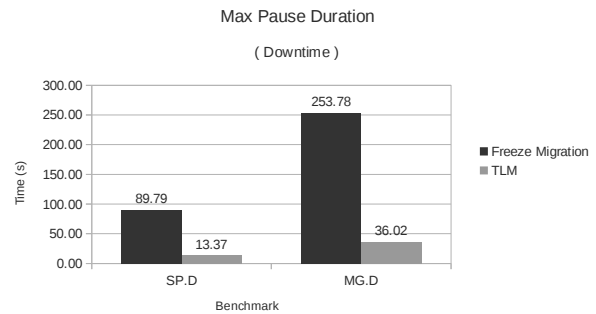
ในการทดลองนี้ ใช้เครื่องคอมพิวเตอร์จำนวน 3 เครื่อง แต่ละเครื่องใช้หน่วยประมวลผลกลางแบบ AMD Opteron Processor 6172 ความเร็ว 2.1 GHz มี 12 แกนประมวลผล (Core) ขนาดหน่วยความจำหลัก 48 GB ขนาดความจุของฮาร์ดดิสก์ 900 GB ระบบปฏิบัติการเป็น Ubuntu 11.10 เคอร์เนล 3.0.0-12-server ติดตั้งโมดูล kvm-kmod-3.0b ทั้ง 3 เครื่องเชื่อมต่อกันด้วยเครือข่ายที่มีความเร็ว 10 Gbps ใช้เก็บข้อมูลฮาร์ดดิสก์ในรูปแบบไฟล์อิมเมจของเวอร์ชวลแมชชีน ด้วยการแชร์ไฟล์ให้กับเครื่องโฮสผ่านโปรโตคอลเอ็นเอฟเอส (NFS) จำนวน 1 เครื่อง และอีก 2 เครื่องใช้งานเป็นโฮสของเวอร์ชวลแมชชีน สร้างเวอร์ชวลแมชชีนที่มีหน่วยประมวลผลกลาง 8 แกนประมวลผล ขนาดหน่วยความจำหลัก 16 GB (สำหรับเบนช์มาร์ก SP.D) และ 36 GB (สำหรับเบนช์มาร์ก MG.D) ขนาดความจุของฮาร์ดดิสก์ 10 GB ระบบปฏิบัติการเกสต์โอเอสคือ Ubuntu 10.04.3 LTS เคอร์เนล 2.6.32-33-generic



รูปที่ 5 โครงสร้างของระบบที่ใช้ในการทดลองไลฟ์ไทม์เกรซัน

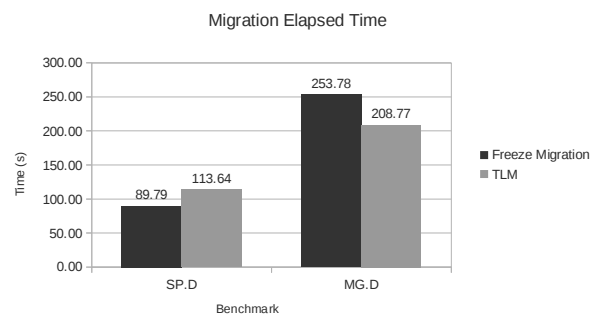
การประเมินประสิทธิภาพทำโดยการเปรียบเทียบเวลาที่ใช้ในการทำงานแบบโอเอ็มพี (OMP) ของโปรแกรม NAS Parallel Benchmark

(NPB) เวอร์ชัน 3.3 เลือกเบนช์มาร์ก SP.D และ MG.D ซึ่งมีการใช้หน่วยความจำขนาด 12.1 GB และ 27.3 GB ตามลำดับ ทำการเปรียบเทียบระหว่างการไลฟ์ไทม์เกรซันแบบเทรค (TLM) ไลฟ์ไทม์เกรซันแบบดั้งเดิม และแบบหยุดแล้วไมเกรซัน



รูปที่ 6 ระยะเวลาที่เวอร์ชวลแมชชีนหยุดการทำงาน

จากผลการทดลองเบื้องต้น ในรูปที่ 6 แสดงระยะเวลาที่เวอร์ชวลแมชชีนหยุดการทำงาน (Downtime) เพื่อทำการคัดลอกข้อมูลไปยังเครื่องโฮสปลายทาง จะเห็นได้ว่าไลฟ์ไทม์เกรซันแบบเทรคใช้เวลาน้อยกว่าแบบหยุดแล้วไมเกรซัน ทั้งสองเบนช์มาร์ก โดยหยุดเวอร์ชวลแมชชีนเป็นเวลา 13.37 วินาที สำหรับเบนช์มาร์ก SP.D และหยุด 36.02 วินาที สำหรับเบนช์มาร์ก MG.D ในขณะที่แบบหยุดแล้วไมเกรซันใช้เวลา 89.79 วินาที และ 253.78 ตามลำดับ ที่เป็นเช่นนี้เพราะที่เอลเอ็มใช้สองเทรคในการส่งข้อมูลไปล่วงหน้า โดยทำงานขนานกัน ทำให้การหยุดการทำงานเพื่อส่งข้อมูลในช่วงสุดท้ายใช้เวลาสั้น

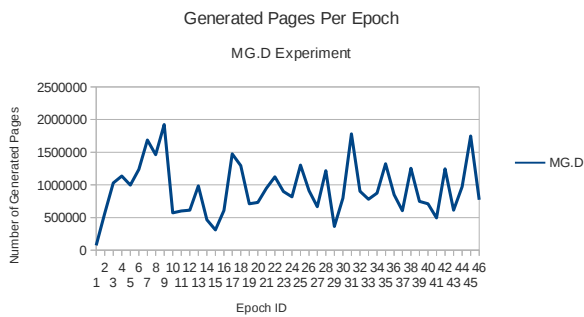


รูปที่ 7 ระยะเวลาทั้งหมดในกระบวนการไมเกรซัน

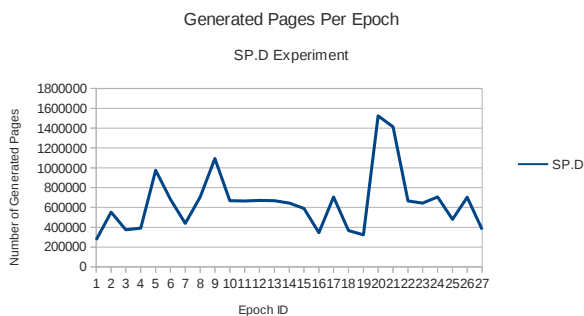
ผลการทดลองในรูปที่ 7 แสดงระยะเวลาตั้งแต่เริ่มต้นไมเกรซันจนจบ (Elapsed Time) จะเห็นได้ว่าในการทดลองที่ใช้เบนช์มาร์ก SP.D ไมเกรซันแบบเดิมใช้เวลา 89.79 วินาที ซึ่งน้อยกว่าไลฟ์ไทม์เกรซันแบบเทรคที่ใช้เวลา 113.64 วินาที แต่ในการทดลองที่ใช้เบนช์มาร์ก MG.D ไมเกรซันแบบเดิมใช้เวลา 253.78 วินาที ซึ่งมากกว่าไลฟ์ไทม์เกรซันแบบเทรคที่ใช้เวลา 208.77 วินาที

ในการทดลองนี้ไคลฟ์ไมเกรชั่นแบบเดิมใช้เวลาานมากเพราะ เบนซ์มาร์กมีการเขียนหน่วยความจำเป็นจำนวนมาก (เป็นแสนเป็นล้านเพจ) อยู่ตลอดเวลา จะทำการไมเกรชั่นสำเร็จก็ต่อเมื่อ เบนซ์มาร์กทำงานจบแล้ว

จากรูปที่ 8 และ 9 เป็นกราฟแสดงปริมาณของเพจ Dirty ที่ถูกสร้างขึ้น จะเห็นได้ว่ากราฟทั้งสองมีความแตกต่างกัน ทั้งนี้ขึ้นอยู่กับแอปพลิเคชันที่นำมาใช้ โดย MG.D มีการสร้างเพจ Dirty เยอะกว่า SP.D มีช่วงของข้อมูลกว้างกว่า และมีจำนวน Epoch มากกว่าด้วย

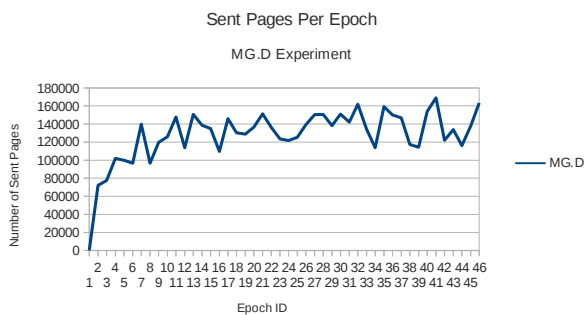


รูปที่ 8 จำนวนเพจที่สร้างขึ้นในแต่ละ Epoch สำหรับ MG.D

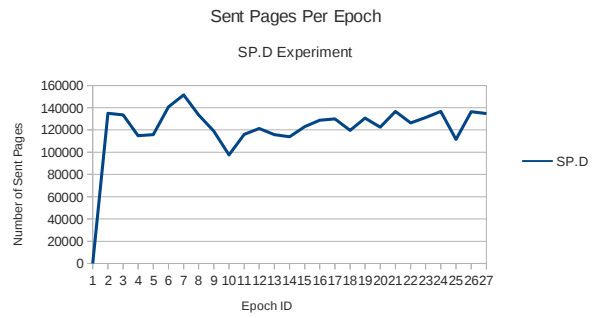


รูปที่ 9 จำนวนเพจที่สร้างขึ้นในแต่ละ Epoch สำหรับ SP.D

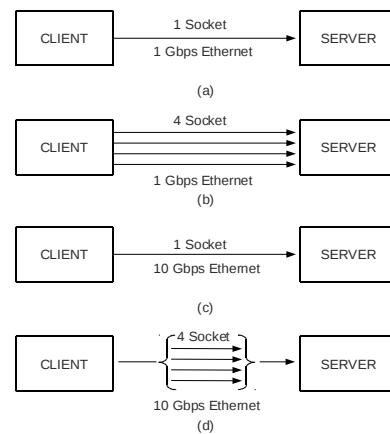
รูปที่ 10 และ 11 แสดงปริมาณของเพจ Dirty ที่เทรคเพจเซฟเวอร์ส่งไปยังเครื่องโฮสปลายทาง โดย MG.D มีการส่งข้อมูลเพจมากกว่า SP.D และมีช่วงของข้อมูลกว้างกว่าด้วย



รูปที่ 10 จำนวนเพจที่เทรคเพจเซฟเวอร์ส่งออกไป สำหรับ MG.D



รูปที่ 11 จำนวนเพจที่เทรคเพจเซฟเวอร์ส่งออกไป สำหรับ MG.D



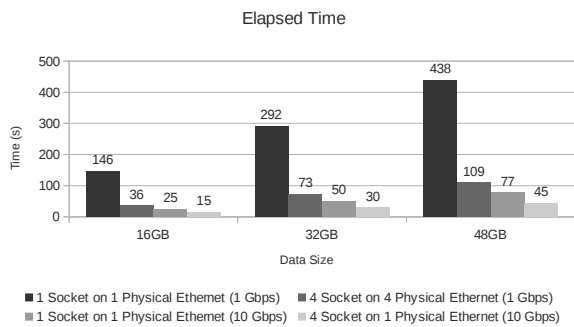
รูปที่ 12 การทดลองเพื่อศึกษาความเป็นไปได้ในการใช้งานไคลบารี

- (a) แบบ 1 ที่ซีพียูซ็อกเก็ต บน 1 ฟิสิคัล 1 Gbps อีเทอร์เน็ต
- (b) แบบ 4 ที่ซีพียูซ็อกเก็ต บน 4 ฟิสิคัล 1 Gbps อีเทอร์เน็ต
- (c) แบบ 1 ที่ซีพียูซ็อกเก็ต บน 1 ฟิสิคัล 10 Gbps อีเทอร์เน็ต
- (d) แบบ 4 ที่ซีพียูซ็อกเก็ต บน 1 ฟิสิคัล 10 Gbps อีเทอร์เน็ต

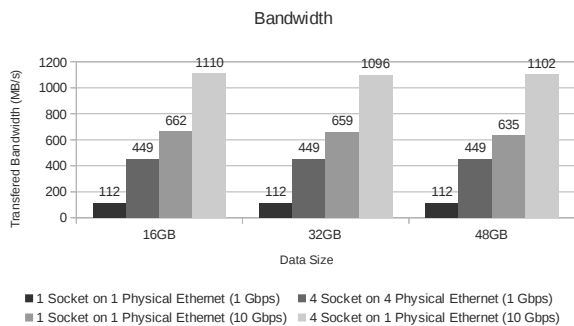
4.2 การทดลองเพื่อศึกษาความเป็นไปได้ในการใช้งานไคลบารี

การทดลองนี้ใช้เครื่องคอมพิวเตอร์จำนวน 2 เครื่อง แต่ละเครื่องใช้หน่วยประมวลผลกลางแบบ AMD Opteron Processor 6172 ความเร็ว 2.1 GHz มี 12 แกนประมวลผล (Core) ขนาดหน่วยความจำหลัก 48 GB ขนาดความจุของฮาร์ดดิสก์ 900 GB ระบบปฏิบัติการเป็น Ubuntu 11.10 เคอร์เนล 3.0.0-12-server เชื่อมต่อกันด้วยเครือข่ายที่มีความเร็ว 1 Gbps และ 10 Gbps เครื่องที่ 1 เป็นไคลเอนต์ เครื่องที่ 2 เป็นเซิร์ฟเวอร์ ในการทดลองนี้ไคลเอนต์มีหน้าที่ส่งข้อมูลขนาดต่าง ๆ ให้กับเซิร์ฟเวอร์ โดยใช้การเชื่อมต่อหลาย ๆ แบบ ขนาดของข้อมูลที่ใช้ในการทดลองได้แก่ 16 GB 32GB และ 48 GB ตามลำดับ การเชื่อมต่อที่ใช้ในการทดลองได้แก่ 1) แบบ 1 ที่ซีพียูซ็อกเก็ต บน 1 ฟิสิคัล 1 Gbps อีเทอร์เน็ต 2) แบบ 4 ที่ซีพียูซ็อกเก็ต บน 4 ฟิสิคัล 1 Gbps อีเทอร์เน็ต 3) แบบ 1 ที่ซีพียูซ็อกเก็ต บน 1 ฟิสิคัล 10 Gbps อีเทอร์เน็ต และ 4) แบบ 4 ที่ซีพียูซ็อกเก็ต บน 1 ฟิสิคัล 10 Gbps อีเทอร์เน็ต การประเมินประสิทธิภาพทำโดยการเปรียบเทียบเวลาที่ไคลเอนต์ใช้ส่งข้อมูลให้กับเซิร์ฟเวอร์ ในแต่ละ

ขนาด และในแต่ละแบบของการเชื่อมต่อ จากนั้นทำการคำนวณและเปรียบเทียบค่าแบนด์วิดท์ จากปริมาณข้อมูลที่ใช้ส่งต่อเวลาที่ใช้ส่ง



รูปที่ 13 เวลาที่โคลนส่งข้อมูลให้กับเซิร์ฟเวอร์



รูปที่ 14 แบนด์วิดท์ระหว่างโคลนและเซิร์ฟเวอร์

จากผลการทดลองพบว่า การส่งข้อมูลผ่าน 10 Gbps Ethernet ที่ใช้การเชื่อมต่อแบบขนานจะมีประสิทธิภาพสูงสุด นั่นคือใช้เวลาในการส่งข้อมูลน้อยที่สุดและมีแบนด์วิดท์สูงสุด โดยใช้เวลา 15 วินาทีในการส่งข้อมูลขนาด 16 GB ใช้เวลา 30 วินาทีในการส่งข้อมูลขนาด 32 GB และใช้เวลา 45 วินาทีในการส่งข้อมูลขนาด 48 GB มีแบนด์วิดท์ประมาณ 1.1 GB/s (8.8 Gbps)

ในกรณีที่ระบบไม่มี 10 Gbps Ethernet สามารถใช้การเชื่อมต่อแบบขนานกับ 1 Gbps Ethernet แทนได้ จากผลการทดลองพบว่า การส่งข้อมูลผ่านการเชื่อมต่อแบบขนาน จำนวน 4 การเชื่อมต่อ ใช้เวลาน้อยกว่า และมีแบนด์วิดท์สูงกว่าการส่งข้อมูลผ่านการเชื่อมต่อแบบ 1 การเชื่อมต่อ 4 เท่า

ผลการทดลองข้างต้น เป็นการช่วยเพิ่มความมั่นใจ ในการที่จะนำไลบรารีการเชื่อมต่อเครือข่ายแบบขนานมาใช้กับการทำไอฟิมเกรชั่น โดยมีแนวโน้มว่า การทำไอฟิมเกรชั่นจะใช้เวลาที่น้อยลงเมื่อใช้วิธีส่งข้อมูลแบบขนาน

6. สรุปผลการทดลองและงานที่จะทำต่อไป

งานวิจัยนี้ทำการปรับปรุงประสิทธิภาพของการทำไอฟิมเกรชั่นให้รองรับเวอร์ชวลแมชชีนที่กำลังรันแอปพลิเคชันที่เน้นการใช้งานหน่วย

ความจำ โดยวิธีไอฟิมเกรชั่นแบบเดิมนั้น ไม่สามารถทำสำเร็จได้ในขณะที่กำลังรันแอปพลิเคชันดังกล่าว อีกวิธีคือหยุดแล้วไอฟิมเกรชั่น วิธีนี้จะต้องหยุดการทำงานของเวอร์ชวลแมชชีนเป็นเวลานาน ผู้วิจัยขอแนะนำวิธีไอฟิมเกรชั่นแบบที่ใช้เทรคมาช่วยในการย้ายข้อมูลบางส่วนไปยังเครื่องปลายทางล่วงหน้า ในขณะที่เดียวกันกับที่เครื่องเวอร์ชวลแมชชีนทำงาน เพื่อให้การหยุดการทำงานของเวอร์ชวลแมชชีน ในขั้นตอนสุดท้ายของการไอฟิมเกรชั่น ใช้เวลาน้อยกว่าวิธีอื่น การประเมินประสิทธิภาพทำโดยใช้โปรแกรม NAS Parallel Benchmarks

ผลการทดลองพบว่าวิธีการทำไอฟิมเกรชั่นแบบเทรค ใช้เวลาหยุดการทำงานของเวอร์ชวลแมชชีนน้อยกว่าวิธีหยุดแล้วไอฟิมเกรชั่น 5 ถึง 7 เท่า

จากผลการทดลองเบื้องต้นเพื่อศึกษาความเป็นไปได้ในการใช้งานไลบรารีพบว่า วิธีการส่งข้อมูลโดยใช้การเชื่อมต่อแบบขนาน จะช่วยให้ส่งข้อมูลได้เร็วขึ้น ใช้เวลาน้อยลง ดังนั้น งานในขั้นต่อไปคือ การนำไลบรารีที่ได้ทำการพัฒนาเพื่อใช้ในการเชื่อมต่อแบบขนาน มาประยุกต์ใช้กับไอฟิมเกรชั่นแบบเทรค โดยมีความคาดหวังว่า การย้ายข้อมูลไปยังเครื่องปลายทางแบบขนานนั้น จะช่วยให้การหยุดการทำงานของเวอร์ชวลแมชชีนใช้เวลาที่น้อยลง และจะช่วยให้ระยะเวลาทั้งหมดในกระบวนการไอฟิมเกรชั่นลดลงด้วย

7. เอกสารอ้างอิง

- [1] <http://aws.amazon.com/ec2/instance-types/>
- [2] Christopher Clark et al. Live migration of virtual machines. In Proceedings of the 2nd conference on Symposium on Networked Systems Design & Implementation - Volume 2 (NSDI'05), Vol. 2. USENIX Association, Berkeley, CA, USA, 273-286.
- [3] vmware, <http://www.vmware.com/files/pdf/VMware-VMotion-DS-EN.pdf>
- [4] <http://www.linux-kvm.org/page/Migration>
- [5] Petter Svärd, Benoit Hudzia, Johan Tordsson, Erik Elmroth, Evaluation of delta compression techniques for efficient live migration of large virtual machines, VEE 2011
- [6] Michael R. Hines, Umesh Deshpande, and Kartik Gopalan. 2009. Post-copy live migration of virtual machines. SIGOPS Oper. Syst. Rev. 43, 3 (July 2009), 14-26.
- [7] Vasinee Siripoonya and Kasidit Chanchio. Thread-Based Live Checkpointing of Virtual Machines. IEEE NCA 2011, pp.155-162