

บทที่ 3

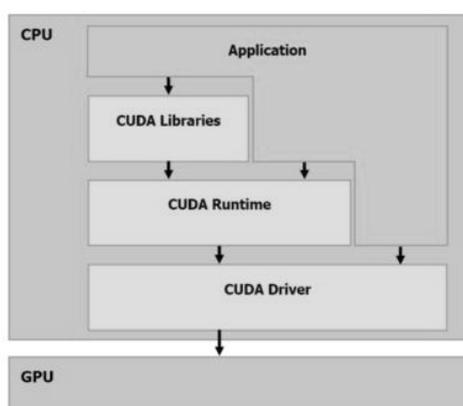
วิธีการดำเนินงานวิจัย

วิทยานิพนธ์ฉบับนี้พัฒนาระบบที่เรียกว่า **เวอร์ชวลคูต้า** มีวัตถุประสงค์เพื่อออกแบบการใช้ทรัพยากรจีพียูร่วมกันจากเครื่องที่ให้บริการทรัพยากรจีพียู (Server) ได้อย่างมีประสิทธิภาพ และวิธีการจัดสรรทรัพยากรจีพียูให้เพียงพอต่อการใช้งานสำหรับคูต้าแอปพลิเคชัน โดยไม่จำกัดว่าผู้ที่ขอรับบริการใช้งานจีพียู (Client) นั้น จะเป็นคูต้าแอปพลิเคชันบนเวอร์ชวลแมชชีนหรือคูต้าแอปพลิเคชันบนเครื่องจริงก็ได้

โดยทั่วไป การเข้าถึงและสั่งงานจีพียูสำหรับการประมวลผลการทำงานของคูต้าแอปพลิเคชันนั้น แอปพลิเคชันจะเรียกใช้งานผ่านคูต้าไดรเวอร์สำหรับการเข้าถึงและสั่งงานไปยังอุปกรณ์จีพียู ไม่ว่าจะเรียกใช้งานในระดับต่ำหรือที่เรียกว่า CUDA driver API คือในส่วนของ CUDA Driver หรือระดับสูงที่เรียกว่า CUDA runtime API คือในส่วนของ CUDA Libraries และ CUDA Runtime ดังภาพที่ 3.1 แต่ทั้งนี้ในระบบคอมพิวเตอร์ที่มีการใช้งานเวอร์ชวลแมชชีนนั้น การเรียกใช้งานผ่านคูต้าไดรเวอร์เพื่อเข้าถึงและสั่งงานไปยังจีพียูจะไม่สามารถเรียกได้จากแอปพลิเคชันที่ทำงานอยู่บนเวอร์ชวลแมชชีน เนื่องจากเกสโอเพอร์เวดิงซิสเต็มของเวอร์ชวลแมชชีนไม่สามารถเข้าถึงจีพียูได้โดยตรง

วิทยานิพนธ์ฉบับนี้จึงได้นำเสนอออกแบบและพัฒนาระบบเวอร์ชวลคูต้าเพื่อสนับสนุนการเข้าถึงและเรียกใช้งานจีพียูขึ้น โดยมีรายละเอียดการออกแบบดังต่อไปนี้

ภาพที่ 3.1 Compute Unified Device Architecture Software Stack

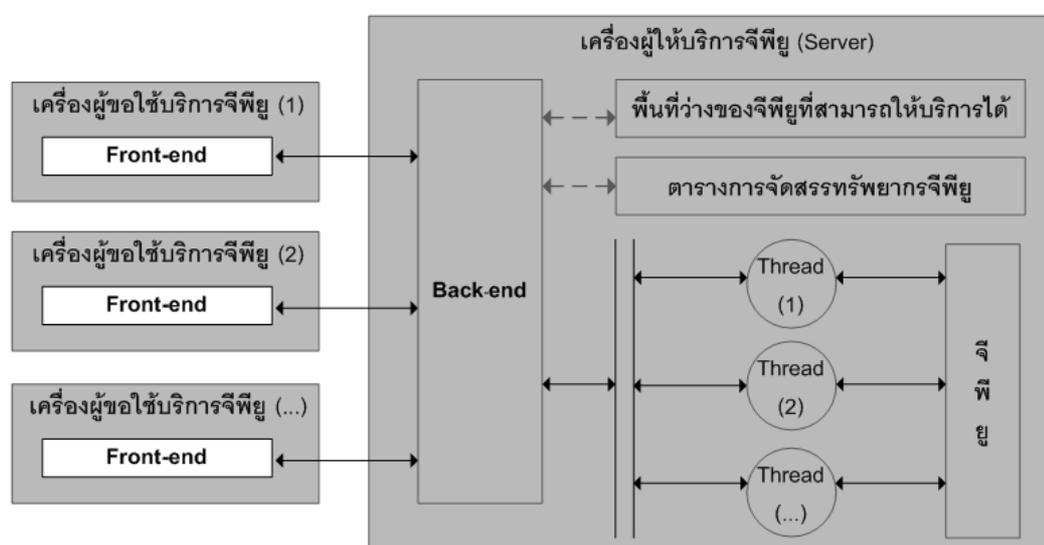


ที่มา : “NVIDIA CUDA Compute Unified Device Architecture” โดย nVidia, 2008, น. 12

3.1 การออกแบบระบบ

วิทยานิพนธ์ฉบับนี้ได้ออกแบบวิธีการเข้าถึงทรัพยากรจ็พ็ญและวิธีการติดต่อสื่อสารระหว่างเครื่องที่ต้องการใช้ทรัพยากรจ็พ็ญกับเครื่องที่ให้บริการทรัพยากรจ็พ็ญโดยได้ออกแบบสถาปัตยกรรมการทำงานของระบบเป็นภาพโดยรวม ดังภาพที่ 3.2

ภาพที่ 3.2 สถาปัตยกรรมภาพรวมของระบบเวอร์ชวลคูด้า



การทำงานของระบบจากภาพที่ 3.2 สามารถอธิบายการลำดับขั้นตอนในการทำงานได้ดังนี้

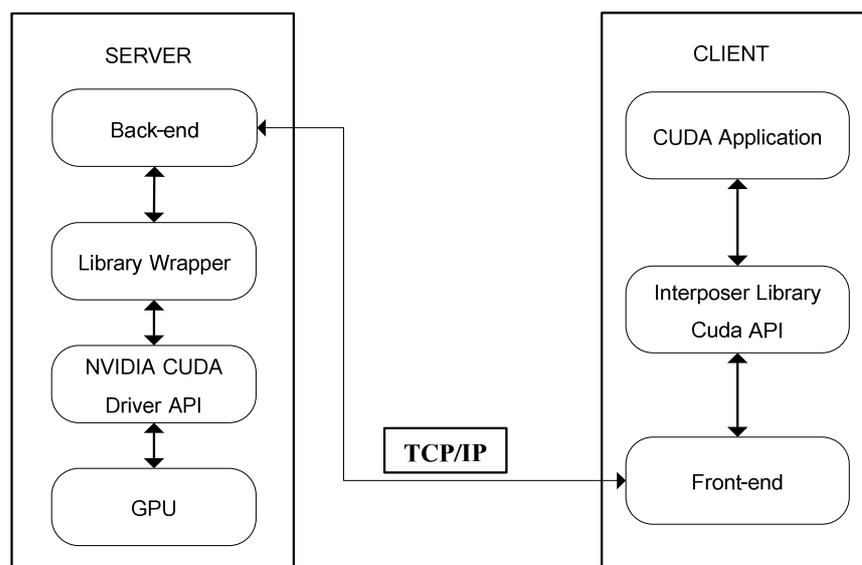
1. คูด้าแอปพลิเคชันบนเครื่องผู้ขอรับบริการจ็พ็ญจะเรียกใช้งานจ็พ็ญผ่านทางพรอนท์เอนด์ไลบรารี (Front-end ดังภาพ) ที่จำลองคูด้าเอพีไอ บนเครื่องของตน
2. เมื่อคำร้องขอใช้จ็พ็ญได้ถูกส่งมาถึงเครื่องผู้ให้บริการจ็พ็ญ เครื่องของผู้ให้บริการจะมีแบคเอนด์ซอฟต์แวร์ (Back-end ดังภาพ) คอยจัดการคำร้องขอการใช้ทรัพยากรจ็พ็ญ
3. แบคเอนด์จะสร้างเทรด (Thread) การทำงานของผู้ขอรับใช้บริการจ็พ็ญโดยจะสร้าง 1 เทรดการทำงานต่อ 1 ผู้ขอรับบริการ เพื่อให้สามารถเรียกใช้งานจ็พ็ญพร้อม ๆ กันได้ครั้งละหลายผู้ขอรับบริการจ็พ็ญ

4. ระหว่างการให้บริการจีพียู แบ็คเอนด์จำเป็นต้องคอยตรวจสอบว่าทรัพยากรจีพียูมีหน่วยความจำสำหรับการให้บริการเพียงพอหรือไม่จากตารางการจัดการทรัพยากรที่แบ็คเอนด์สร้างขึ้น

3.1.1 การทำงานของผู้ให้บริการทรัพยากรจีพียู (Server) และผู้ใช้บริการทรัพยากรจีพียู (Client)

จากการอธิบายการทำงานภาพรวมสถาปัตยกรรมการทำงานของระบบเวอร์ชวลคูด้าสามารถจำแนกการทำงานของระบบเวอร์ชวลคูด้าออกเป็น 2 ส่วนการทำงานดังแสดงในภาพที่ 3.3 คือ 1) ส่วนการทำงานของผู้ให้บริการทรัพยากรจีพียู (Server) ซึ่งมีส่วนประกอบในการทำงานของ แบ็คเอนด์ (Back-end) และ Library Wrapper 2) ส่วนการทำงานของผู้ใช้ทรัพยากรจีพียู (Client) มีส่วนประกอบในการทำงานของ Interposer Library CUDA API และ ฟรอนท์เอนด์ (Front-end) โดยมีรายละเอียดในการทำงานของ Server และ Client แสดงในภาพที่ 3.3 ดังต่อไปนี้

ภาพที่ 3.3 การทำงานระหว่างเครื่องผู้ใช้บริการจีพียู (Client) กับเครื่องที่ให้บริการจีพียู (Server)



3.1.1.1 การทำงานของผู้ให้บริการทรัพยากรจีพียู (Server) มีรายละเอียดหน้าที่ในการทำงานดังนี้

1. แบคเอนด์ (Back-end) ทำหน้าที่เป็นตัวกลางผสมงานและจัดการ การเข้าถึงจีพียูของแต่ละคำสั่งตามคูด้าแอปพลิเคชันที่ถูกส่งมาจากฟรอนท์เอนด์ของเครื่องผู้ขอรับบริการจีพียู โดยนำคำสั่งที่ได้รับไปประมวลผลและเรียกใช้งานจีพียูให้ประมวล หลังจากประมวลผลเสร็จสิ้น แบคเอนด์จะนำผลที่ได้ส่งกลับไปยังฟรอนท์เอนด์ของเครื่องที่ขอใช้บริการจีพียูของแต่ละเครื่อง

สำหรับการเรียกใช้งานคูด้าแอปพลิเคชันของแบคเอนด์นั้นจำเป็นต้องมีการสร้างแอปพลิเคชันสำหรับเรียกใช้งานไว้ที่ผู้ให้บริการจีพียูสำหรับการเรียกใช้งานบนจีพียู โดยมีรายละเอียดดังต่อไปนี้

ก. เริ่มจากนำโปรแกรมสำหรับประมวลผลการปฏิบัติงานบนจีพียูแล้วนำมาแปลโปรแกรม (Compile) ด้วย NVCC เพื่อให้ได้ Cubin File โดยใช้ตัวเลือกให้ออกมาเป็น PTX Code

ข. นำ Cubin File ดังกล่าวไปวางไว้บนเครื่องของผู้ให้บริการทรัพยากรจีพียู สำหรับแบคเอนด์เรียกใช้งานโดยผ่าน Library Wrapper

2. Library Wrapper ทำหน้าที่แปลงคำสั่งจากฟรอนท์เอนด์ที่ต้องการเรียกใช้งานจีพียูด้วยคำสั่งของ CUDA Driver API ผ่าน NVIDIA CUDA Library ที่ถูกติดตั้งไว้บนเครื่องผู้ให้บริการจีพียูให้อยู่ในรูปของคำสั่ง CUDA Driver API ตามปกติ สำหรับใช้ในการติดต่อสื่อสารกับอุปกรณ์จีพียู เพื่อผสมผสานการทำงานและประมวลผลการทำงานตามแต่ละคำสั่งที่ได้รับจากฟรอนท์เอนด์

3.1.1.2 การทำงานของผู้ขอใช้ทรัพยากรจีพียู (Client) มีรายละเอียดในการทำงานดังนี้

1. การทำงานของผู้ขอใช้ทรัพยากรจีพียูเริ่มต้นเมื่อมีการเรียกใช้งานคูด้าแอปพลิเคชันโดยส่วนที่เป็นคูด้าโค้ดที่ต้องการเข้าถึงจีพียูเพื่อใช้จีพียูในการประมวลผลการทำงาน ซึ่งโดยปกติแล้วคำสั่ง CUDA Driver API จะถูกเรียกใช้งานผ่าน NVIDIA CUDA Library แต่ในกรณีที่ผู้ขอใช้ทรัพยากรจีพียูไม่มีอุปกรณ์จีพียูเป็นของตนเองจึงจำเป็นต้องเข้าถึงจีพียู โดยผ่าน Interposer Library CUDA API สำหรับเรียกใช้งานคำสั่ง CUDA Driver API ที่อยู่บนเครื่องของผู้ให้บริการทรัพยากรจีพียู

2. Interposer Library CUDA API ทำหน้าที่เป็นตัวกลางผสมงานจัดการ การเรียกใช้งาน NVIDIA CUDA Library จากคู่มือแอปพลิเคชันของผู้ให้บริการจีพียู โดยไลบรารีที่ทำหน้าที่เป็นตัวกลางนี้จะถูกนำไปวางไว้บนเครื่องของผู้ให้บริการจีพียู และเมื่อไลบรารีถูกเรียกใช้งานจะทำการเตรียมข้อมูลที่ต้องส่งไปให้แบ็คเอนด์ เช่น ตัวแปร ขนาดของข้อมูล และประเภทคำสั่งของคู่มือ นำทุกอย่างมารวมกันเป็นแพคเกจแล้วจึงส่งแพคเกจนั้น ๆ ต่อให้ฟรอนท์เอนด์ทำงานต่อไป

3. ฟรอนท์เอนด์ (Front-end) ทำหน้าที่จัดการ การติดต่อสื่อสารระหว่างเครื่องผู้ให้บริการจีพียูกับเครื่องผู้ให้บริการจีพียู โดยเริ่มจากการสร้างช่องทางการติดต่อด้วยซ็อกเก็ต (Socket) ผ่านโปรโตคอลที่ซีพี/ไอพี (TCP/IP) แล้วจึงจัดส่งแพคเกจที่ได้รับจาก Interposer Library CUDA API ส่งไปตามช่องทางที่ถูกสร้างขึ้นไปที่แบ็คเอนด์ของเครื่องผู้ให้บริการทรัพยากรจีพียู ซึ่งการเลือกใช้โปรโตคอลที่ซีพีสำหรับฟรอนท์เอนด์นั้น เนื่องจากที่ซีพีมีกลไกในการตรวจสอบการสูญหายของข้อมูลระหว่างการรับส่งข้อมูลต่างจากยูดีพี (UDP) ทำให้เชื่อมั่นได้ว่าข้อมูลที่ผ่านการรับส่งผ่านที่ซีพีจะไม่เกิดการสูญหายระหว่างทาง

สำหรับฟังก์ชันการทำงานของ CUDA Driver API ที่อยู่ใน Library Wrapper และ Interposer Library CUDA API งานวิจัยนี้ได้รวบรวมฟังก์ชันพื้นฐานที่สำคัญของ CUDA Driver API ที่นำมาใช้สำหรับการประมวลผลปฏิบัติงานของคู่มือแอปพลิเคชัน โดยส่วนที่เป็นฟังก์ชันสำหรับประมวลผลปฏิบัติงานบนจีพียูนั้นจะถูกแปลโปรแกรม (Compile) ให้อยู่ในรูปแบบของ Cubin File และวางอยู่ในส่วนของผู้ให้บริการทรัพยากรจีพียูสำหรับรอการถูกเรียกใช้งานจากแบ็คเอนด์ ฟังก์ชันพื้นฐานที่สำคัญของ CUDA Driver API ที่นำมาใช้ได้แก่

- | | |
|---------------------------------|-----------------------|
| ➤ cuInit | ➤ cuParamSetv |
| ➤ cuDeviceGet | ➤ cuParamSetSize |
| ➤ cuCtxCreate | ➤ cuFuncSetBlockShape |
| ➤ cuModuleLoad | ➤ cuLaunchGrid |
| ➤ cuModuleGetFunction | ➤ cuCtxDetach |
| ➤ cuMemAlloc | ➤ cuMemFree |
| ➤ cuMemcpyHtoD/
cuMemcpyDtoH | |

โดยมีเหตุผลในการเลือกใช้ฟังก์ชันพื้นฐานที่สำคัญของ CUDA Driver API ทั้ง 14 ฟังก์ชันนี้ เนื่องจากเป็นฟังก์ชันที่สำคัญพื้นฐานสำหรับการสั่งงานให้จีพียูสามารถประมวลผลการทำงานได้ และเวลาทำงานวิจัยที่มีอยู่อย่างจำกัด

3.1.2 การจัดการสรรกรการเข้าถึงทรัพยากรจีพียู

ส่วนนี้กล่าวถึงการออกแบบการจัดการสรรกรการเข้าถึงจีพียูของผู้ให้บริการจีพียูโดยแสดงขั้นตอนการทำงานโดยรวมดังภาพที่ 3.4 , 3.5 และ 3.6 ซึ่งการออกแบบการจัดการสรรกรทรัพยากรจีพียูของระบบเวอร์ชวลคูด้า ได้ออกแบบให้มีตารางสำหรับการจัดการสรรกรทรัพยากรเพื่อหลีกเลี่ยงการเกิดปัญหาการติดตาย อันเนื่องมาจากเหตุการณ์ที่แบ็คเอนด์ทำการจองพื้นที่บนจีพียูทันทีหลังจากได้รับคำสั่งจากฟรอนท์เอนด์ และคำสั่งที่รับมาพร้อม ๆ กันจากหลายเวอร์ชวลแมชชีนโดยไม่ใช้คำสั่งสำหรับสั่งให้จีพียูประมวลผลซึ่งอาจจะทำให้เกิดการรอคำสั่งประเภทดังกล่าวจนเกิดการติดตาย หรือเกิดการยกเลิกการประมวลผลของคูด้าแอปพลิเคชันกรณีมีพื้นที่ไม่เพียงพอต่อการประมวลผล อีกทั้งระบบเวอร์ชวลคูด้ายังได้ให้ความสำคัญกับลำดับในการเข้าถึงทรัพยากรจีพียูจากการร้องขอใช้ทรัพยากร และคำนึงถึงการเกิดปัญหาการติดตาย (Deadlock)

วิทยานิพนธ์ฉบับนี้จึงได้ออกแบบอัลกอริทึมสำหรับจัดการสรรกรการเข้าถึงทรัพยากรจีพียูออกเป็น 3 ขั้นตอน คือ 1) การจัดการสรรกรทรัพยากรให้เพียงพอต่อการเรียกใช้งาน 2) การจัดสิทธิในการเข้าถึงทรัพยากรจากลำดับการขอเข้าใช้งานจีพียู 3) ตรวจสอบคำสั่งการปฏิบัติงาน cuLaunchGrid กรณีเรียกใช้งานซ้ำ โดยอธิบายรายละเอียดแต่ละขั้นตอนดังนี้

3.1.2.1 การจัดการสรรกรทรัพยากรให้เพียงพอต่อการเรียกใช้งาน

การจัดการสรรกรทรัพยากรจีพียูสำหรับผู้ขอใช้บริการทรัพยากรจีพียูนั้น วิทยานิพนธ์นี้ได้นำอัลกอริทึมของนายธนาคาร (Banker's Algorithm) เข้ามาประยุกต์ในการจัดการสรรกรทรัพยากรเพื่อใช้ในการตรวจสอบทรัพยากรว่ามีเพียงพอต่อการให้บริการหรือไม่ อีกทั้งยังสามารถหลีกเลี่ยงปัญหาการติดตาย (Deadlock) ซึ่งการเลือกอัลกอริทึมของนายธนาคารมาใช้เนื่องจาก อัลกอริทึมของนายธนาคาร เป็นอัลกอริทึมที่ใช้หลีกเลี่ยงการเกิดปัญหาการติดตายแต่ไม่ใช่อัลกอริทึมเพื่อแก้ปัญหาการติดตายเหมือนกับ Detection Algorithm จึงเหมาะสมกับระบบเวอร์ชวลคูด้าที่จัดทำขึ้น โดยมีขั้นตอนการจัดการสรรกรทรัพยากรดังนี้

1. เมื่อแบ็คเอนด์รับคำสั่งคูด้าแอปพลิเคชันมาจากเครื่องผู้ขอใช้บริการแบ็คเอนด์ จะทำหน้าที่ตรวจสอบก่อนว่าคำสั่งคูด้านั้นมีหน้าที่ใด กรณีเป็นคำสั่งเพื่อจองพื้นที่บนจีพียู

(cuMemAlloc) แบริคเอ็นด์จะนำตัวแปร ขนาดของตัวแปร และหมายเลขเครื่องผู้ให้บริการจีพียู ไปเก็บไว้ยังตารางการจัดสรรทรัพยากร (ดังตารางที่ 3.1) เพื่อรอและนำไปใช้งานต่อไป อีกทั้งนำ ข้อมูลที่ถูกส่งมาจากพรอนท์เอนด์ด้วยคำสั่ง cuMemcpyHtoD พักไว้บนแบริคเอ็นด์สำหรับรอ เรียกใช้งานในการถ่ายโอนข้อมูลไปยังจีพียูเมื่อมีการสั่งให้ปฏิบัติงานด้วยคำสั่ง cuLaunchGrid เกิดขึ้นจากการสั่งงานของพรอนท์เอนด์ การนำทุกคำสั่งคู่ด้าแอปพลิเคชันจากพรอนท์เอนด์มาเก็บ ในตารางการจัดสรรทรัพยากร เนื่องจากการทำงานด้วยวิธีนี้สามารถป้องกันการเกิดปัญหาติดตาย และแอปพลิเคชันที่ประมวลผลไม่สำเร็จเนื่องจากพื้นที่สำหรับการใช้ประมวลผลไม่เพียงพอต่อการ ใช้งาน

2. กรณีที่แบริคเอ็นด์รับคำสั่งคู่ด้าแอปพลิเคชัน ซึ่งส่งมาจากเครื่องผู้ให้บริการ และเป็นคำสั่งเพื่อเรียกใช้งานในฟังก์ชันของคู่ด้าแอปพลิเคชัน ซึ่งคำสั่งดังกล่าวคือคำสั่ง cuLaunchGrid ในส่วนการทำงานของแบริคเอ็นด์จะทำการตรวจสอบก่อนว่าตัวแปรที่ฟังก์ชันต้อง นำไปใช้นั้นมีอะไรบ้างและเมื่อทราบแล้ว แบริคเอ็นด์จะรวมพื้นที่ทั้งหมดของตัวแปรดังกล่าวและ นำไปตรวจสอบกับพื้นที่ของจีพียูที่ยังเหลืออยู่ว่ามีขนาดเพียงพอกับเนื้อที่ทั้งหมดของตัวแปรที่จะ นำไปใช้หรือไม่ หากเพียงพอแบริคเอ็นด์จะจัดการเรียกการทำงานผ่าน NVIDIA CUDA Library เพื่อจองพื้นที่ให้กับตัวแปรทั้งหมดพร้อมกันทำสัญลักษณ์ให้กับตัวแปรในตารางการจัดสรร ทรัพยากรว่าตัวแปรดังกล่าวได้ถูกจองพื้นที่ไว้แล้วในหน่วยความจำของจีพียู การทำสัญลักษณ์ว่า มีการจองพื้นที่ให้กับตัวแปรไปแล้วในตารางการจัดสรรทรัพยากรนี้มีวัตถุประสงค์เพื่อ เมื่อมีการ เรียกใช้งานตัวแปรเดิมซ้ำแบริคเอ็นด์จะไม่นำตัวแปรดังกล่าวไปจองพื้นที่ในหน่วยความจำจีพียูซึ่ง ตัวแปรทุกตัวจะถูกยกเลิกออกจากตารางการจัดสรรทรัพยากรและหน่วยความจำจีพียูด้วยคำสั่ง cuMemFree

3.1.2.2 การจัดสิทธิในการเข้าถึงทรัพยากรจากลำดับการขอเข้าใช้งานจีพียู

ส่วนนี้กล่าวถึงการออกแบบการจัดการการเข้าถึงจีพียูของผู้ให้บริการจีพียู โดยการ จัดสิทธิสำหรับการเข้าถึงจีพียูจากลำดับการขอเข้าใช้งานจีพียูซึ่งการจัดการดังกล่าวนี้ จะจัดการ ในระดับคำสั่งของคู่ด้า ตามการเรียกใช้ของผู้ให้บริการจีพียู ด้วยวิธีมาก่อนจะได้รับสิทธิในการ ประมวลผลก่อน (FIFO) ซึ่งการจัดลำดับความสำคัญนี้จะให้ความสำคัญที่คำสั่ง cuLaunchGrid โดยมีรายละเอียดตามขั้นตอนดังนี้

1. เมื่อผู้ให้บริการจีพียูได้รับคำสั่ง cuLaunchGrid ในส่วนการทำงานของ แบริคเอ็นด์จะทำหน้าที่ให้ลำดับความสำคัญในการประมวลผลและหมายเลขเครื่องผู้ให้บริการจีพียู

2. หลังจากให้ลำดับความสำคัญในการประมวลผลแล้ว แบ็คเอนด์จะทำการสำรวจคำสั่งการให้ปฏิบัติงานของ cuLaunchGrid จากตารางการจัดสรรทรัพยากรและสั่งให้ cuLaunchGrid ที่มีลำดับการปฏิบัติงานสูงสุดปฏิบัติงาน โดยต้องผ่านการตรวจสอบใน 3.1.2.1) เป็นที่เรียบร้อยแล้ว

3. เมื่อการประมวลผลตามคำสั่งการปฏิบัติงานของ cuLaunchGrid เสร็จสิ้นหรือแบ็คเอนด์ได้รับคำสั่งให้ยกเลิกการจองพื้นที่ของตัวแปรในจีพียู แบ็คเอนด์จำเป็นต้องสำรวจตารางการจัดสรรทรัพยากรโดยทำตาม ขั้นตอนที่ 2) สำหรับการจัดสิทธิในการเข้าถึงทรัพยากรจากลำดับการขอเข้าใช้งานจีพียู จนกระทั่งไม่มีคำสั่ง cuLaunchGrid เหลืออยู่ในตารางจัดสรรทรัพยากร

3.1.2.3 ตรวจสอบคำสั่งการปฏิบัติงาน cuLaunchGrid กรณีเรียกใช้งานซ้ำ

ส่วนนี้กล่าวถึงการตรวจสอบคำสั่งปฏิบัติงาน cuLaunchGrid กรณีเรียกใช้งานซ้ำจากผู้ขอรับบริการจีพียู ซึ่งส่งงานผ่านมาทางฟรอนท์เอนด์ และเมื่อคำสั่งดังกล่าวมาถึงแบ็คเอนด์ แบ็คเอนด์มีหน้าที่ต้องตรวจสอบ ตามรายละเอียดดังนี้

1. ตรวจสอบข้อมูลตัวแปรแต่ละตัวที่คำสั่ง cuLaunchGrid ต้องการเรียกใช้งาน เพื่อนำไปประมวลผลว่ามีการจองพื้นที่แล้วบนจีพียูแล้วหรือไม่ กรณีที่มีการจองพื้นที่แล้วแบ็คเอนด์จะไม่ทำการจองพื้นที่ซ้ำอีกครั้ง

2. ตรวจสอบข้อมูลตัวแปรที่เป็น Input แต่ละตัวที่คำสั่ง cuLaunchGrid ต้องการเรียกใช้งานว่ามีการคัดลอกข้อมูลไปยังจีพียูแล้วหรือไม่ กรณีที่มีการคัดลอกแล้ว แบ็คเอนด์จะไม่ทำการคัดลอกซ้ำลงไปอีก

เมื่อทำการตรวจสอบตามข้อ 1. และ 2. เสร็จเรียบร้อยแล้วและผลปรากฏว่าเป็นไปตามเงื่อนไข คือตัวแปรแต่ละตัวได้ถูกจองพื้นที่บนจีพียูแล้ว และตัวแปรที่เป็น Input ได้ถูกคัดลอกข้อมูลไปยังจีพียูแล้ว นั้นแสดงว่า cuLaunchGrid เคยถูกเรียกใช้งานมาก่อนและข้อมูลต่าง ๆ ที่อยู่บนจีพียูยังไม่เคยถูกล้างออกไป ซึ่งจะเป็นเหตุผลให้แบ็คเอนด์เรียกใช้งาน cuLaunchGrid ได้ทันที โดยที่ไม่ต้องไปเริ่มทำคำสั่ง cuInit ถึง cuMemcpyHtoD ใหม่อีกครั้ง จึงเป็นการลดเวลาในการทำงานในขั้นตอนดังกล่าวสำหรับแบ็คเอนด์ หรือแม้แต่ฟรอนท์เอนด์หากต้องการเรียกใช้งาน cuLaunchGrid เดิมซ้ำ ก็ไม่จำเป็นต้องสั่งงาน cuInit ถึง cuMemcpyHtoD มียังแบ็คเอนด์อีก ซึ่งจะเป็นการลดเวลาในการถ่ายโอนข้อมูลมายังแบ็คเอนด์เป็นอย่างมาก

3.1.2.4 ตัวอย่างการจัดสรรการเข้าถึงทรัพยากรจีพียู

ตารางที่ 3.1 ตัวอย่างตารางการจัดสรรทรัพยากรตามคำสั่งคูด้าแอฟพลีเคชัน

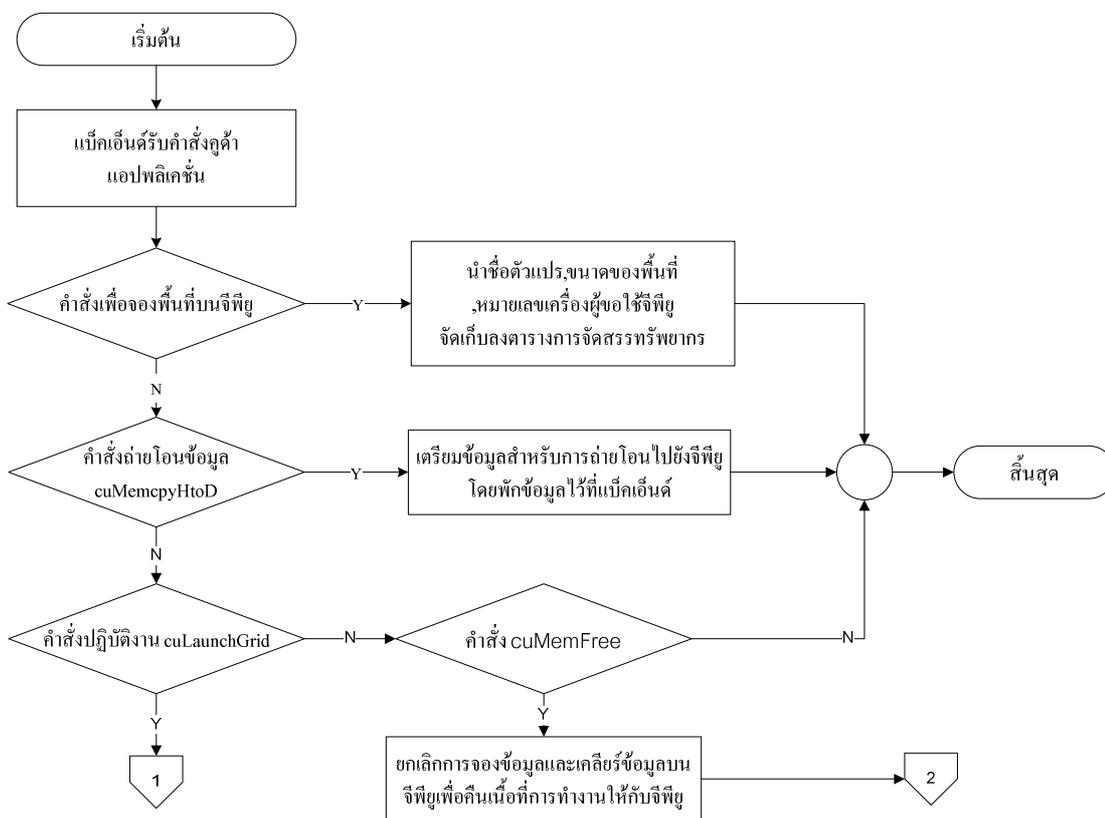
คำสั่ง	ชื่อตัวแปร	ขนาดข้อมูล	เครื่องผู้ขอใช้ จีพียู	ลำดับการ ทำงาน	พื้นที่ถูก จอง
cuMemcpyHtoD	A	20MB	1	-	ไม่
cuMemcpyHtoD	A	30MB	2	-	ใช่
cuMemcpyHtoD	B	20MB	2	-	ใช่
cuMemcpyHtoD	B	50MB	1	-	ไม่
cuLaunchGrid	-	0	1	1	-
cuLaunchGrid	-	0	2	2	-

จากตารางที่ 3.1 สามารถอธิบายตัวอย่างการทำงานจัดสรรการเข้าถึงทรัพยากรจีพียูได้ดังนี้ คือตารางที่ 3.1 เป็นตารางการจัดสรรทรัพยากรจีพียูที่ฝังอยู่บนแคเคเอ็นดีของเครื่องผู้ให้บริการจีพียู โดยมีเหตุการณ์สมมุติดังนี้ เมื่อมีการเรียกใช้งานคำสั่งคูด้าฟังก์ชันในการจองพื้นที่สำหรับตัวแปร A จากเครื่องหมายเลข 1 ขนาดของข้อมูล 20 MB และมีประเภทการจองพื้นที่เป็น cuMemcpyHtoD คือขอจองพื้นที่เพื่อคัดลอกข้อมูลจากจีพียูไปยังจีพียู จากพรอนท์เอนด์ของเครื่องเวอร์ชวลแมชีน แคเคเอ็นดีจะทำการตรวจสอบคำสั่งที่ส่งมาก่อนว่าเป็นคำสั่งปฏิบัติงาน cuLaunchGrid หรือไม่ กรณีตรวจแล้วไม่ใช่จะนำคำสั่ง cuMemcpyHtoD ชื่อตัวแปร A หมายเลขเครื่องและขนาดของข้อมูลไปเก็บไว้ในตารางการจัดสรร ดังแสดงในตารางที่ 3.1 และทำแบบเดียวกันนี้ไปเรื่อย ๆ จนกว่าคำสั่งที่ส่งมาจากพรอนท์เอนด์จะเป็นคำสั่ง cuLaunchGrid ซึ่งในตารางที่ 3.1 ได้เก็บไว้ 4 คำคือ A,B ของเครื่องที่ 1 และ 2 กรณีตรวจสอบแล้วเป็นคำสั่ง cuLaunchGrid จริง แคเคเอ็นดีจะยังไม่ทำการรันคูด้าแอฟพลีเคชันแต่จะตรวจสอบพื้นที่ก่อนว่ามีเพียงพอสำหรับการประมวลผลหรือไม่ โดยจะรวมค่าตัวแปรทั้งหมดที่ต้องใช้สำหรับ cuLaunchGrid นั้น ๆ กรณีมีพื้นที่เพียงพอแคเคเอ็นดีจะทำการประมวลผลคูด้าแอฟพลีเคชันนั้น ๆ ทันที แต่ถ้าหากมีพื้นที่ไม่เพียงพอจะนำคำสั่ง cuLaunchGrid ไปเก็บไว้ในตารางการจัดสรรทรัพยากรจีพียู และให้ลำดับในการรอประมวลผลจนกว่ามีพื้นที่ว่างพอให้ประมวลผล จากตัวอย่างกำหนดให้พื้นที่ในการประมวลผลจีพียูเหลือเพียง 50 MB เพราะฉะนั้นจะเห็นได้ว่าเมื่อมีคำสั่ง

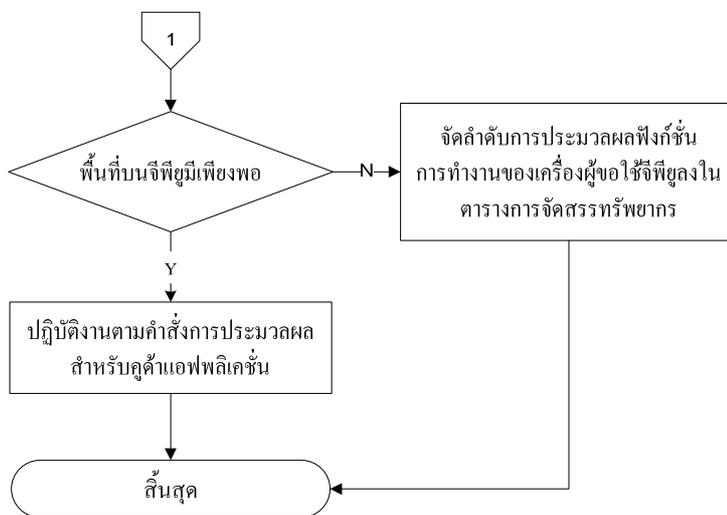
cuLaunchGrid ที่ส่งมาจากพรอนท์เอนด์ จากตารางที่ 3.1 แถวที่ 5 ซึ่งเป็นของเครื่องหมายเลขที่ 1 พบว่าเมื่อแบคเอนด์ทำการตรวจสอบพื้นที่ที่ต้องใช้ในการประมวลผลสำหรับเครื่องหมายเลข 1 นั้น พื้นที่ของตัวแปร A บวกกับพื้นที่ของ ตัวแปร B สำหรับใช้ในการประมวลผลคู่ค่าแอฟลิเคชัน สำหรับคำสั่ง cuLaunchGrid นั้นมีพื้นที่มากเกินไปที่จีพียูจะสามารถให้บริการได้ แบคเอนด์จึงนำคำสั่ง cuLaunchGrid ของเครื่องหมายเลข 1 ไปเก็บไว้ในตารางการจัดสรรทรัพยากรจีพียู พร้อมกับบรรทัดของสถิติที่ได้รับสำหรับการทำงานเมื่อจีพียูมีพื้นที่เพียงพอ และต่อมาแบคเอนด์ได้รับคำสั่ง cuLaunchGrid ของเครื่องหมายเลข 2 ปรากฏว่าจากการตรวจสอบพบว่า มีพื้นที่เหลือเพียงพอในการประมวลผลจึงปล่อยให้ cuLaunchGrid ของเครื่องหมายเลข 2 ทำงาน และเมื่อมีพื้นที่ว่างจากการที่คู่ค่าแอฟลิเคชันก่อนหน้าประมวลผลเสร็จ cuLaunchGrid ของเครื่องหมายเลข 1 จะถูกทำงานทันทีเช่นกัน (ดังภาพที่ 3.4 , 3.5 และ 3.6)

จากตารางที่ 3.1 สามารถอธิบายตัวอย่างตรวจสอบคำสั่งการปฏิบัติงาน cuLaunchGrid กรณีเรียกใช้งานซ้ำได้ดังนี้ คือ กรณีที่เครื่องผู้ขอใช้จีพียูหมายเลข 2 เข้าไปสั่งงานจีพียูตามคำสั่ง cuLaunchGrid แสดงว่ามีพื้นที่บนจีพียูเพียงพอสำหรับให้บริการ ทำให้ตัวแปร A และ B สามารถจองพื้นที่ และคัดลอกข้อมูลลงไปที่จีพียูได้ ขณะที่ตัวแปร A และ B จองพื้นที่ และคัดลอกข้อมูลลงไปที่จีพียูได้นั้น แบคเอนด์จะทำสัญลักษณ์การจองพื้นที่ไว้บนตารางการจัดสรรว่าตัวแปรดังกล่าวเคยจองพื้นที่ และคัดลอกข้อมูลลงไปที่จีพียูแล้ว โดยจะถูกยกเลิกข้อมูลก็ต่อเมื่อพบคำสั่ง cuMemFree ทั้งนี้ทำให้เมื่อเครื่องผู้ขอใช้จีพียูหมายเลข 2 สั่ง cuLaunchGrid อีกครั้ง แบคเอนด์จะไปตรวจสอบข้อมูลตัวแปร A และ B ที่ตารางการจัดสรร ซึ่งจะพบว่าข้อมูลทั้งสองตัวแปรเคยถูกจองพื้นที่และมีข้อมูลอยู่บนจีพียูแล้ว ซึ่งจะทำให้แบคเอนด์สามารถสั่ง cuLaunchGrid ได้ทันทีโดยไม่ต้องจองพื้นที่และคัดลอกข้อมูลลงไปที่จีพียูซ้ำอีกครั้ง (ดังภาพที่ 3.7)

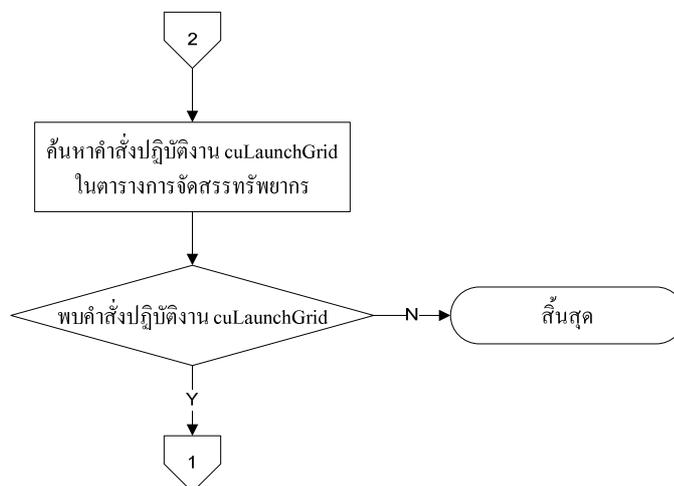
ภาพที่ 3.4 ภาพรวมการจัดสรรการเข้าถึงทรัพยากรจีพียูของระบบเวอร์ชวลคูด้า สำหรับขั้นตอนการตรวจสอบแต่ละประเภทคำสั่ง



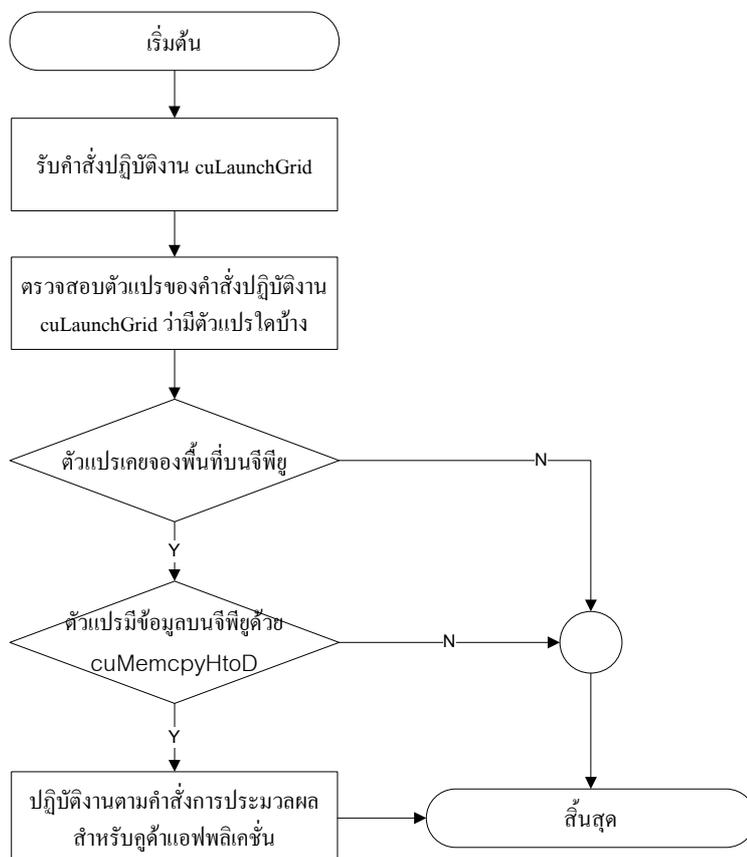
ภาพที่ 3.5 ภาพรวมการจัดสรรการเข้าถึงทรัพยากรจีพียูของระบบเวอร์ชวลคูด้า สำหรับขั้นตอนการตรวจสอบพื้นที่บนจีพียูเพื่อประมวลผล



ภาพที่ 3.6 ภาพรวมการจัดสรรการเข้าถึงทรัพยากรจีพียูของระบบเวอร์ชวลคูด้า สำหรับขั้นตอน
การตรวจสอบคำสั่งปฏิบัติงานในตารางจัดสรร



ภาพที่ 3.7 ภาพตรวจสอบคำสั่งการปฏิบัติงาน cuLaunchGrid กรณีเรียกใช้งานซ้ำ
ของระบบเวอร์ชวลคูด้า



3.1.3 ตัวอย่างโปรแกรมสำหรับเรียกใช้งานจีพียูด้วยคูด้าแอปพลิเคชันโดยผ่านระบบเวอร์ชวลคูด้า

ระบบเวอร์ชวลคูด้าที่ได้พัฒนาขึ้นมา นั้น จะทำงานโดยเข้าถึงจีพียูด้วยคำสั่งของคูด้าแบบ CUDA driver API โดยจะใช้ฟังก์ชันพื้นฐานที่ได้กล่าวไว้ในข้อ 3.1.1.2 โดยแสดงการเรียกใช้งาน ดังนี้

```
#include "functioncuda_client.h"

void matrixMul(int argc, char** argv)
{
    int sockfd=connectServer();

    // set seed for rand()
    srand(2006);

    unsigned int uiWA, uiHA, uiWB, uiHB, uiWC, uiHC;
    int iSizeMultiple = 1;
    uiWA = WA * iSizeMultiple;
    uiHA = HA * iSizeMultiple;
    uiWB = WB * iSizeMultiple;
    uiHB = HB * iSizeMultiple;
    uiWC = WC * iSizeMultiple;
    uiHC = HC * iSizeMultiple;

    unsigned int size_A = uiWA * uiHA;
    unsigned int mem_size_A = sizeof(float) * size_A;
    float* h_A = (float*)malloc(mem_size_A);
    unsigned int size_B = uiWB * uiHB;
    unsigned int mem_size_B = sizeof(float) * size_B;
    float* h_B = (float*)malloc(mem_size_B);

    // initialize host memory
    randomInit(h_A, size_A);
```

```
randomInit(h_B, size_B);

cuInitialize(sockfd);
cuLoadFunction("matrixMul",sockfd);

unsigned int size_C = uiWC * uiHC;
unsigned int mem_size_C = sizeof(float) * size_C;

cudaMalloc((void**) &d_C, mem_size_C);*/
cuMemAllocDevice(mem_size_C,"d_C",sockfd);
cuMemAllocDevice(mem_size_A,"d_A",sockfd);
cuMemAllocDevice(mem_size_B,"d_B",sockfd);

cuMemcpy("d_A",h_A,mem_size_A,HtoD,sockfd);
cuMemcpy("d_B",h_B,mem_size_B,HtoD,sockfd);

cudaCallFunction(sockfd,BLOCK_SIZE,BLOCK_SIZE,1,GridX,GridY);
cuMemcpy("d_C",h_C,mem_size_C,DtoH,sockfd));

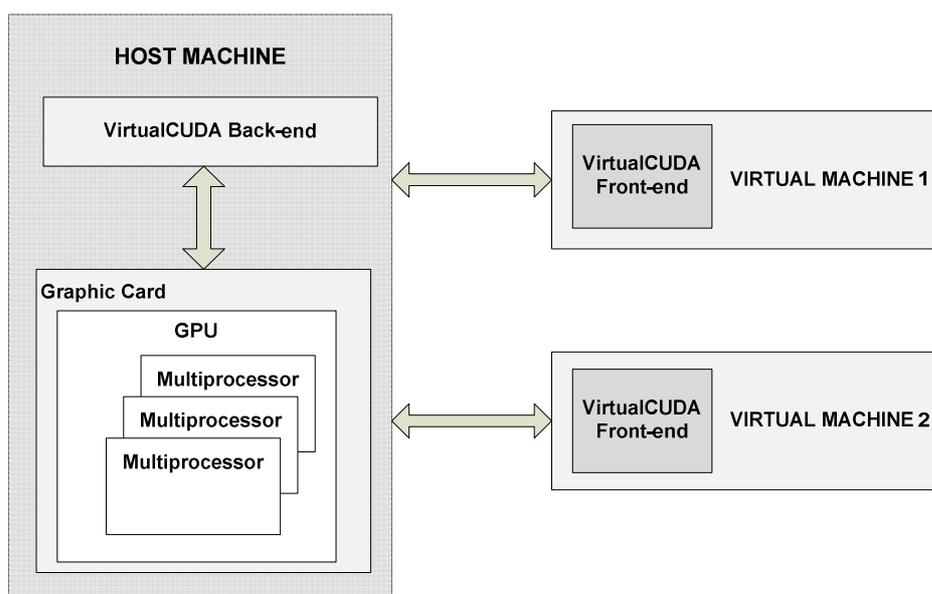
free(h_A);
free(h_B);
free(h_C);
cudaMemFree("d_A",sockfd);
cudaMemFree("d_B",sockfd);
cudaMemFree("d_C",sockfd);
}
```

3.2 วิธีการทดลอง

3.2.1 เครื่องมือที่ใช้ในการทดลอง

ในการทดลองนี้เรากำหนดให้มีเวอร์ชวลแมชชีนสองเครื่องสามารถเรียกใช้งานทรัพยากรจํานวนหนึ่งของโฮสแมชชีน โดยผ่านระบบเวอร์ชวลคูด้า ดังภาพที่ 3.8 และสองเวอร์ชวลแมชชีน ดังกล่าวจะถูกรันอยู่บนเครื่องโฮสเครื่องเดียวกัน ซึ่งแสดงรายละเอียดของ โฮสแมชชีน เวอร์ชวลแมชชีน และ Graphic Card ดังต่อไปนี้

ภาพที่ 3.8 การทำงานของโฮสแมชชีนและสองเวอร์ชวลแมชชีน



1. โฮสแมชชีนสำหรับแบ็คเอนด์ มีรายละเอียดดังนี้
 - Architecture : x86_64
 - Processor : Intel(R) Xeon(R) CPU E5530 @ 2.40GHz
 - CPU : 4
 - Memory capacity : 8GB
 - OS : Ubuntu 10.10
 - Kernel : 2.6.35-24-server
2. เวอร์ชวลแมชชีนแบบ KVM เครื่องที่ 1 มีรายละเอียดดังนี้

- Architecture : x86_64
 - Processor : QEMU Virtual CPU version 0.12.5
 - Memory capacity : 1GB
 - OS : Ubuntu 9.04
 - Kernel: 2.6.28-13-generic
3. เวอร์ชวลแมชชีนแบบ KVM เครื่องที่ 2 มีรายละเอียดดังนี้
- Architecture : x86_64
 - Processor : QEMU Virtual CPU version 0.12.5
 - Memory capacity : 1GB
 - OS : Ubuntu 9.04
 - Kernel : 2.6.28-13-generic
4. ไฮเปอร์ไวเซอร์ kvm-88
5. Graphic Card
- รุ่น : NVIDIA GeForce 8800GTS
 - CUDA Driver Version : 3.2
 - CUDA Cores : 96
 - Core Clock (MHz) : 500
 - Shader Clock (MHz) : 1200
 - Memory Clock (MHz) : 800
 - Memory Amount : 640MB
 - Memory Interface : 320-bit
 - Memory Bandwidth (GB/sec) : 64

3.2.2 การออกแบบการทดลอง

ในส่วนนี้อธิบายถึงการออกแบบการทดลองสำหรับการขอรับบริการใช้งานทรัพยากรจีพียูของเวอร์ชวลแมชชีนโดยแบ่งการทดลองออกเป็น 3 แบบคือ 1) การทดลองประมวลผลผ่านซีพียูของเวอร์ชวลแมชชีนเปรียบเทียบกับ การประมวลผลผ่านเวอร์ชวลคูด้า 2) การทดลองเพื่อวัด

ระยะเวลาที่เรียกใช้งานจีพียูผ่านระบบเวอร์ชวลคูด้า 3) การทดลองการใช้งานจีพียูร่วมกันของเวอร์ชวลแมชชีนผ่านระบบเวอร์ชวลคูด้า โดยมีรายละเอียดในการทดลองแต่ละแบบดังต่อไปนี้

3.2.2.1 แบบที่ 1 การทดลองประมวลผลผ่านจีพียูของเวอร์ชวลแมชชีนเปรียบเทียบกับ การประมวลผลผ่านเวอร์ชวลคูด้า มีรายละเอียดดังต่อไปนี้

• วิธีการทดลอง

1. ทำการรันแอปพลิเคชันบนเครื่องเวอร์ชวลแมชชีนโดยที่แอปพลิเคชันดังกล่าวประมวลผลบนจีพียู จากนั้นทำการจับเวลาทั้งหมดที่ใช้ไปสำหรับการประมวลผล
2. นำเอาแอปพลิเคชันเดียวกันที่เคยประมวลผลบนจีพียูมาดัดแปลงให้เป็นคูด้าแอปพลิเคชันเพื่อให้สามารถประมวลผลบนจีพียูได้ โดยการประมวลผลของคูด้าแอปพลิเคชันดังกล่าวจำเป็นต้องประมวลผลผ่านระบบเวอร์ชวลคูด้า จากนั้นทำการจับเวลาทั้งหมดที่ใช้ไปสำหรับการประมวลผล
3. นำเวลาที่ได้จากการประมวลผลทั้งในข้อ 1 และข้อ 2 มาเปรียบเทียบกัน

• เป้าหมายการทดลอง

1. เพื่อวัดประสิทธิภาพในการประมวลผลและการใช้ประโยชน์ (Utilization) ที่เพิ่มมากขึ้น เมื่อนำจีพียูเข้ามาช่วยในการประมวลผล
2. เพื่อแสดงให้เห็นความสามารถในการทำงานสำหรับการเข้าถึงจีพียูของระบบเวอร์ชวลคูด้าจากเวอร์ชวลแมชชีน

• สมมติฐานการทดลอง

การประมวลผลผ่านคูด้าแอปพลิเคชันเพื่อเรียกใช้งานจีพียูผ่านเวอร์ชวลคูด้า การประมวลผลมีประสิทธิภาพที่ดีกว่าการประมวลผลผ่านจีพียูของเวอร์ชวลแมชชีน

3.2.2.2 แบบที่ 2 การทดลองเพื่อวัดระยะเวลาที่เรียกใช้งานจีพียูผ่านระบบเวอร์ชวลคูด้า มีรายละเอียดดังต่อไปนี้

• วิธีการทดลอง

1. รันคูด้าแอปพลิเคชันสำหรับเรียกใช้งานจีพียูบนเวอร์ชวลแมชชีนผ่านระบบเวอร์ชวลคูด้า

2. วัดระยะเวลาการทำงานของแบ็คเอนด์ และฟรอนท์เอนด์ ตั้งแต่ Initialization (คำสั่ง cuInit) ถึง Function Call (คำสั่ง cuLaunchGrid) จากนั้นทำการเปรียบเทียบระยะเวลาที่ได้ดังกล่าวระหว่างแบ็คเอนด์ กับฟรอนท์เอนด์

3. วัดระยะเวลาการทำงานของฟังก์ชันการคัดลอกข้อมูลแบบ Host to Device (จากซีพียูไปซีพียู) ของแบ็คเอนด์ และฟรอนท์เอนด์ จากนั้นทำการเปรียบเทียบระยะเวลาที่ได้ดังกล่าวระหว่างแบ็คเอนด์ กับฟรอนท์เอนด์

4. วัดระยะเวลาการทำงานของฟังก์ชันการคัดลอกข้อมูลแบบ Device to Host (จากซีพียูไปซีพียู) ระหว่างแบ็คเอนด์ และฟรอนท์เอนด์ จากนั้นทำการเปรียบเทียบระยะเวลาที่ได้ดังกล่าวระหว่างแบ็คเอนด์ กับฟรอนท์เอนด์

• เป้าหมายการทดลอง

1. เพื่อวัดผลและวิเคราะห์โอเวอร์เฮดการประมวลผลของคูด้าแอปพลิเคชันสำหรับแบ็คเอนด์ และฟรอนท์เอนด์กรณีเรียกใช้งานจีพียูผ่านระบบ เวอร์ชวลคูด้า

2. เพื่อวัดผลและวิเคราะห์โอเวอร์เฮดการถ่ายโอนข้อมูลระหว่างแบ็คเอนด์ และฟรอนท์เอนด์ กรณีเรียกใช้งานจีพียูผ่านระบบเวอร์ชวลคูด้า

• สมมติฐานการทดลอง

1. ระยะเวลาการประมวลผลของแบ็คเอนด์ให้ผลที่ดีกว่าการประมวลผลของฟรอนท์เอนด์

2. การถ่ายโอนข้อมูลไปมาระหว่างแบ็คเอนด์ และฟรอนท์เอนด์มีผลทำให้ระยะเวลาในการประมวลผลของคูด้าแอปพลิเคชันบนเวอร์ชวลแมชชีนเพิ่มมากขึ้นตามขนาดของข้อมูลที่ถูกถ่ายโอน

3.2.2.3 แบบที่ 3 การทดลองการใช้งานจีพียูร่วมกันของเวอร์ชวลแมชชีนผ่านระบบเวอร์ชวลคูด้า มีรายละเอียดดังต่อไปนี้

• วิธีการทดลอง

1. ทำการรันคูด้าแอปพลิเคชันจากเครื่องเวอร์ชวลแมชชีนผ่านระบบเวอร์ชวลคูด้า เพื่อร้องขอใช้บริการจีพียูมายังเครื่องผู้ให้บริการ กรณีเวอร์ชวลแมชชีนมากกว่า 1 เครื่อง แต่แบ็คเอนด์สามารถขอใช้บริการจีพียูได้เพียงครั้งละ 1 เวอร์ชวลแมชชีน โดยที่แบ็คเอนด์ไม่แบ่งเทรดการทำงานให้กับเวอร์ชวลแมชชีนแต่ละเครื่องที่เข้ามาขอใช้บริการจีพียูทำให้เวอร์ชวลแมชชีนต้องต่อคิวเพื่อเข้าใช้งานจีพียู

2. เพิ่มประสิทธิภาพแก่ระบบเวอร์ชวลคู่มือในการให้บริการจีพียู สำหรับการรันคู่มือแอปพลิเคชันบนเครื่องเวอร์ชวลแมชชีน ด้วยการเพิ่มความสามารถให้ กับแบ็คเอนด์ในฝั่งของผู้ให้บริการจีพียูให้สามารถรองรับการทำงานของเวอร์ชวลแมชชีน มากกว่า 1 เครื่อง กรณีขอเข้าใช้งานจีพียูพร้อมกันเพื่อประโยชน์ในการใช้ทรัพยากรจีพียูแบบเต็มประสิทธิภาพมากที่สุด

3. กำหนดให้คู่มือแอปพลิเคชันที่ใช้ในการประมวลผลทั้งในข้อ 1 และข้อ 2 เป็นคู่มือแอปพลิเคชันเดียวกัน

• เป้าหมายการทดลอง

1. เพื่อวัดประสิทธิภาพและเปรียบเทียบการทำงานในข้อ 1 และข้อ 2 ว่ามีการใช้ประโยชน์ (Utilization) สำหรับการเรียกใช้งานจีพียูกรณีใช้งานร่วมกันเพิ่มขึ้นหรือไม่

2. เพื่อวัดผลและวิเคราะห์โอเวอร์เฮดสำหรับเวลาที่ใช้ไปในการประมวลผลของคู่มือแอปพลิเคชัน กรณีที่ขอใช้ทรัพยากรจีพียูร่วม

• สมมติฐานการทดลอง

1. เวอร์ชวลคู่มือสามารถรองรับเวอร์ชวลแมชชีนกรณีเรียกใช้งานจีพียูพร้อมกันสองเครื่องได้

2. การเข้าใช้งานจีพียูพร้อมกันของสองเวอร์ชวลแมชชีนได้โดยผ่านเวอร์ชวลคู่มือมีประสิทธิภาพในการทำงานที่ดีกว่าการเข้าใช้งานจีพียูโดยผ่านเวอร์ชวลคู่มือครั้งละหนึ่งเวอร์ชวลแมชชีน

3.2.3 การวิเคราะห์ข้อมูลและการวัดผล

1. วิเคราะห์ประสิทธิภาพเปรียบเทียบระหว่างการประมวลผลบนจีพียูและบนจีพียูผ่านเวอร์ชวลคู่มือของเวอร์ชวลแมชชีน

2. วัดผลและวิเคราะห์โอเวอร์เฮดของการใช้งานคู่มือแอปพลิเคชัน สำหรับแบ็คเอนด์ และฟรอนท์เอนด์

3. วิเคราะห์ประสิทธิภาพเปรียบเทียบระหว่างการใช้งานจีพียูที่แบ็คเอนด์อนุญาตให้เข้าถึงได้เพียงครั้งละ 1 เวอร์ชวลแมชชีนกับการใช้งานจีพียูพร้อมกันครั้งละมากกว่า 1 เวอร์ชวลแมชชีน

4. วัดผลและวิเคราะห์โอเวอร์เฮดระหว่างการใช้งานจีพียูที่ 1 เวอร์ชวลแมชชีนกับการใช้งานจีพียูพร้อมกัน 2 เวอร์ชวลแมชชีน