

## Multilevel Cache Design Considerations

- ❑ Design considerations for L1 and L2 caches are very different
  - Primary cache should focus on **minimizing hit time** in support of a shorter clock cycle
    - Smaller with smaller block sizes
  - Secondary cache(s) should focus on **reducing miss rate** to reduce the penalty of long main memory access times
    - Larger with larger block sizes
- ❑ The miss penalty of the L1 cache is significantly reduced by the presence of an L2 cache – so it can be smaller (i.e., faster) but have a higher miss rate
- ❑ For the L2 cache, hit time is less important than miss rate
  - The L2\$ hit time determines L1\$'s miss penalty
  - L2\$ local miss rate >> than the global miss rate

CSE431 L20&21 Improving Cache Performance.17

Irwin, PSU, 2005

เครื่องคอมพิวเตอร์ยุคใหม่จะใช้งานแคช (CACHE) หลายระดับเพื่อเพิ่มประสิทธิภาพการทำงานของระบบโดยรวม ในการใช้งานนั้น โปรเซสเซอร์ (PROCESSOR) ใช้งานข้อมูลที่อยู่ในแคชปฐมภูมิ (PRIMARY CACHE – L1 CACHE) เท่านั้น ไม่อาจใช้งานข้อมูลที่อยู่ในระดับต่ำลงไปได้โดยตรง ดังนั้นแคชในระดับนี้ควรจะมีค่าหน่วยเวลาในการใช้งานต่ำเพื่อที่โปรเซสเซอร์ไม่ต้องเสียเวลาในการรอข้อมูลนาน ดังนั้นการออกแบบแคชปฐมภูมินี้จึงมีความต้องการให้มีการเข้าถึงข้อมูลด้วยเวลาที่น้อยที่สุด นั่นคือแคชปฐมภูมิจะต้องใช้หน่วยความจำที่มีระยะเวลาในการเข้าถึงข้อมูล (ACCESS TIME) ต่ำๆ หรือเป็นหน่วยความจำที่สามารถใช้กับคาบสัญญาณเวลาต่างๆ ได้ เนื่องจากหน่วยความจำประเภทนี้มักจะมีราคาสูงทำให้ไม่เหมาะที่จะใช้แคชขนาดใหญ่ ในกรณีที่ข้อมูลที่ต้องการใช้งานไม่อยู่บนแคชปฐมภูมิ เรียกว่าเกิดแคชมิสส์ (CACHE MISS) ระบบจะต้องนำข้อมูลที่ต้องการใช้งานมาจากหน่วยความจำระดับต่ำลงมาซึ่งในกรณีนี้คือแคชทุติยภูมิ (SECONDARY CACHE – L2 CACHE) มาบรรจุในแคชปฐมภูมิ โปรเซสเซอร์จึงจะสามารถใช้งานได้ ในช่วงที่นำข้อมูลจากแคชทุติยภูมิมาไว้ที่แคชปฐมภูมินั้น โปรเซสเซอร์จะต้องหยุดการทำงานเพื่อรอข้อมูล ดังนั้นการออกแบบแคชในระดับปฐมภูมินี้จึงไม่นิยมที่จะให้มีขนาดของบล็อก (BLOCK) ใหญ่มากๆ เพื่อให้เวลาในการโอนถ่ายข้อมูลต่ำที่สุด และโปรเซสเซอร์สามารถดำเนินงานต่อไปได้

ในกรณีที่ข้อมูลที่โปรเซสเซอร์ต้องการไม่อยู่ในแคชปฐมภูมิ จะเกิดการใช้งานข้อมูลในแคชทุติยภูมิขึ้น ถ้าหากว่าข้อมูลที่โปรเซสเซอร์ต้องการอยู่ในแคชทุติยภูมิ เวลาที่ต้องใช้เมื่อเกิดเหตุการณ์นี้เรียกว่ามิสส์เพนัลตี (MISS PENALTY) ของแคชปฐมภูมิ คือเวลาที่ต้องใช้ในการโอนถ่ายข้อมูลจากแคชทุติยภูมิไปยังแคชปฐมภูมิ ซึ่งเท่ากับเวลาที่ใช้ในการเข้าถึงข้อมูลของแคชทุติยภูมิซึ่งเป็นเวลาที่น้อยกว่าเวลาในการเข้าถึงข้อมูลของหน่วยความจำหลัก ถ้าหากว่าข้อมูลที่โปรเซสเซอร์ต้องการไม่อยู่ในแคชปฐมภูมิและแคชทุติยภูมิ มิสส์เพนัลตีจะเป็นระยะเวลาที่นานกว่ามาก ซึ่งเท่ากับเวลาที่ใช้ในการโอนถ่ายข้อมูลจากหน่วยความจำหลักมายังแคชทุติยภูมิก่อนที่เวลาที่ใช้ในการโอนถ่ายข้อมูลจากแคชทุติยภูมิไปยังแคชปฐมภูมิ

การออกแบบใช้งานแคชปฐมภูมิและแคชทุติยภูมิในระบบที่เป็นแคชหลายระดับนั้นมียุทธศาสตร์ของแคชในแต่ละระดับที่แตกต่างกัน สำหรับโครงสร้างที่เป็นแคชสองระดับนี้ จุดประสงค์ของการออกแบบแคชปฐมภูมิคือลดเวลาการเข้าถึงข้อมูลของโปรเซสเซอร์ให้ใช้เวลาหรือจำนวนสัญญาณเวลาน้อยที่สุด ในขณะที่จุดประสงค์ของการออกแบบแคชทุติยภูมิคือลดอัตราการมิสส์ (MISS RATE) เพื่อลดเวลาในการโอนถ่ายข้อมูล

เมื่อพิจารณาถึงปฏิสัมพันธ์ระหว่างแคชทั้งสองระดับจะพบว่ามิสส์เพนัลตีของแคชปฐมภูมิลดลงอย่างมากเมื่อมีการใช้แคชทุติยภูมิในระบบ แคชปฐมภูมิจึงมีขนาดเล็กและมีอัตราการมิสส์ที่สูงขึ้นได้ เมื่อพิจารณาที่แคชทุติยภูมิพบว่าเวลาในการเข้าถึงข้อมูลเป็นประเด็นที่สำคัญน้อยลงเพราะมีแคชปฐมภูมิอยู่ เพราะเวลาในการเข้าถึงข้อมูลของแคชทุติยภูมิมีผลกระทบต่อมิสส์เพนัลตีของแคชปฐมภูมิ ไม่ได้มีผลกระทบโดยตรงกับเวลาในการโอนถ่ายข้อมูลไปยังโปรเซสเซอร์โดยตรง

การใช้งานระบบที่ใช้แคชหลายระดับนั้นมีความซับซ้อนอยู่บ้าง ประการหนึ่งคือมีการมิสส์และอัตราการมิสส์อยู่หลายประเภท เช่นมีอัตราการมิสส์ของแคชปฐมภูมิหรืออัตราการมิสส์แบบครอบคลุม (GLOBAL MISS RATE) ซึ่งก็คือสัดส่วนการอ้างถึงข้อมูลที่มีมิสส์ (ไม่พบข้อมูลในแคช) จากแคชทุติยภูมิ มีอัตราการมิสส์ของแคชทุติยภูมิซึ่งก็คือสัดส่วนของมิสส์ในแคชทุติยภูมิหารด้วยจำนวนการเข้าถึงข้อมูล อัตราการมิสส์นี้เรียกว่าอัตราการมิสส์แบบเฉพาะที่ (LOCAL MISS RATE) ของแคชทุติยภูมิ อัตราการมิสส์แบบเฉพาะที่ของแคชทุติยภูมิมิสูงกว่าอัตราการมิสส์แบบครอบคลุมเป็นอย่างมาก เพราะว่าแคชปฐมภูมิเป็นด่านรับการเข้าถึงข้อมูลไปแล้วเป็นอันมากโดยเฉพาะอย่างยิ่งกับการเข้าถึงข้อมูลที่อยู่ในบริเวณใกล้เคียงกันหรือการเข้าถึงข้อมูลที่ถูกใช้ซ้ำบ่อยๆ จากตัวอย่างที่นำเสนอในห้วงเรียนอัตราการมิสส์แบบเฉพาะที่จะอยู่ที่  $0.5\% / 2.0\% = 25\%$  ซึ่งเป็นอัตราที่สูงมาก แต่โดยภาพรวมแล้วระบบไม่ได้มีผลกระทบมากนักเพราะอัตราการมิสส์แบบครอบคลุมเป็นตัวหลักที่บอกให้ทราบความบ่อยครั้งในการเข้าถึงข้อมูลของระบบ