
6.004 Handbook

Written and Illustrated by Margaret Chong
Edited by Joseph Wong

Please feel free to email me at mchong@mit.edu if you have any comments or suggestions *for this handbook*. Please contact your TA if you have any conceptual questions.

Acknowledgements

Thank you Professor Terman for inspiring me to pursue a career in circuit design and systems engineering.

Thank you Professor Ward for giving me the opportunity to TA – something I've wanted to do since I was a freshman.

Thank you Jeffrey Sheldon, Brandon Kam, and James Sun for proofreading this handbook.

Table of Contents

Binary	Intro to Stacks
Octal and Hexadecimal	Stack Macros
Binary Addition	Stack Frames
Binary Subtraction	Two-level Stack example
Information	The Memory Hierarchy
Huffman Encoding	Temporal Locality
Hamming Distance	Fully-Associative Cache
The Digital Abstraction	Direct-Mapped Cache
Combinational Devices	N-Way Set-Associative Cache
Static Discipline	Spatial Locality and Data Block Size
NFETs	Replacement Policies
PFETs	Write Policies
CMOS Inverter	Virtual Memory
Transistor Level CMOS Logic	Linear Page Table
CMOS Logic Gates	Hierarchical Page Map
Boolean Algebra	Translation Look-aside Buffer
Building Circuits with Logic Gates	Caches – Virtual or physical addresses
Contamination Delay (coming eventually)	MMU Overview
Propagation Delay (coming eventually)	
Sum of Products	
Karnaugh Maps	
Lenient Combinational Logic	
Mux	
ROM/PLA	
Latch	
Flip Flop	
Flip Flop Circuits	
Finite State Machines (FSMs)	
Metastability	
Unpipelined Circuits	
Pipelined Circuits	

Binary

There are two types of binary used in 6.004 – unsigned and signed (also known as two's complement). Let's take a binary number **1100₂** and figure out its decimal equivalent.

Unsigned $2^3 \ 2^2 \ 2^1 \ 2^0 \Rightarrow 2^3 + 2^2 = 8 + 4 = 12$

1	1	0	0
---	---	---	---

Signed $2^3 \ 2^2 \ 2^1 \ 2^0 \Rightarrow (-1)(2^3) + 2^2 = -8 + 4 = -4$

- (1) The space to the left of the decimal point is for 2^0 , increasing by a power of two as you go further left, decreasing by a power of two as you go right
- (2) For each 1 in the binary number, add the corresponding power of two to the sum
- (3) Unsigned binary - everything in the sum is positive
- (4) Signed binary – the leftmost space in the number is negative

Note: If the leftmost space is a 0, then everything in the sum is positive. For example,

$$1100_2 = -8 + 4 = -4$$

$$01100_2 = 8 + 4 = 12 \quad \leftarrow \text{the leftmost 0 means we don't count the } -16$$

Octal and Hexadecimal

Octal is base 8 and is made of groups of 3 bits

Hexadecimal is base 16 and is made of groups of 4 bits

Binary

1	1	1	1	0	0	1	1	0	1	1	1
7			4			6			7		
F				3				7			

Octal

Hexadecimal

Useful Conversion Tables:

Binary	Octal
000	0
001	1
010	2
011	3
100	4
101	5
110	6
111	7

Binary	Hex	Binary	Hex
0000	0	1000	8
0001	1	1001	9
0010	2	1010	A
0011	3	1011	B
0100	4	1100	C
0101	5	1101	D
0110	6	1110	E
0111	7	1111	F

Signed Binary Addition

Binary addition is the same as regular addition, except you only have 0's and 1's.

$$\begin{array}{r}
 3 \\
 + 5 \\
 \hline
 8
 \end{array}
 \Rightarrow
 \begin{array}{r}
 1 \\
 00011 \\
 + 00101 \\
 \hline
 0
 \end{array}
 \Rightarrow
 \begin{array}{r}
 11 \\
 00011 \\
 + 00101 \\
 \hline
 00
 \end{array}
 \Rightarrow
 \begin{array}{r}
 111 \\
 00011 \\
 + 00101 \\
 \hline
 000
 \end{array}
 \Rightarrow
 \begin{array}{r}
 111 \\
 00011 \\
 + 00101 \\
 \hline
 1000
 \end{array}
 \Rightarrow
 \begin{array}{r}
 111 \\
 00011 \\
 + 00101 \\
 \hline
 01000
 \end{array}$$

\uparrow
 $1+1 = 2 = 10$ in binary, so carry the 1

Most machines have a fixed number of bits. In the example above, the operands are both 5 bits wide, therefore the result can only be expressed in 5 bits. This introduces the possibility that the sum might be incorrect.

$$\begin{array}{r}
 3 \\
 + 5 \\
 \hline
 ??
 \end{array}
 \Rightarrow
 \begin{array}{r}
 1 \\
 0011 \\
 + 0101 \\
 \hline
 0
 \end{array}
 \Rightarrow
 \begin{array}{r}
 11 \\
 0011 \\
 + 0101 \\
 \hline
 00
 \end{array}
 \Rightarrow
 \begin{array}{r}
 111 \\
 0011 \\
 + 0101 \\
 \hline
 000
 \end{array}
 \Rightarrow
 \begin{array}{r}
 111 \\
 0011 \\
 + 0101 \\
 \hline
 1000
 \end{array}$$

-8 is incorrect!

We cannot express 8 in four bits

$$\begin{array}{r}
 -8 \\
 + -8 \\
 \hline
 ??
 \end{array}
 \Rightarrow
 \begin{array}{r}
 1000 \\
 + 1000 \\
 \hline
 0000
 \end{array}
 \Rightarrow
 \begin{array}{r}
 1 \\
 1000 \\
 + 1000 \\
 \hline
 0000
 \end{array}$$

0 is incorrect!

We cannot express -16 in four bits

Signed Binary Subtraction

Let's do this by example, say we're subtracting 3 minus 5

- (1) Represent 3 and 5 in binary
- (2) Figure out the binary representation for *negative* 5 by flipping the bits and adding one
- (3) Compute 3 plus negative 5

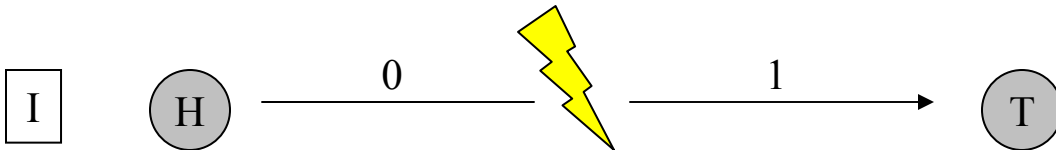
$$\begin{array}{r}
 3 \\
 - 5 \\
 \hline
 -2
 \end{array}
 \xrightarrow{1}
 \begin{array}{r}
 3 = 00011 \\
 5 = 00101
 \end{array}
 \xrightarrow{2}
 \begin{array}{r}
 3 = 00011 \\
 -5 = 11010 \\
 + 00001 \\
 \hline
 11011
 \end{array}
 \xrightarrow{3}
 \begin{array}{r}
 11 \\
 00011 \\
 + 11011 \\
 \hline
 11110
 \end{array}$$

Note that, like in binary addition, we are constrained to a fixed number of bits. Therefore, you might encounter cases where the "flip the bits and add one" cannot be represented in some fixed number of bits. You may also find that the result cannot be represented in some fixed number of bits.

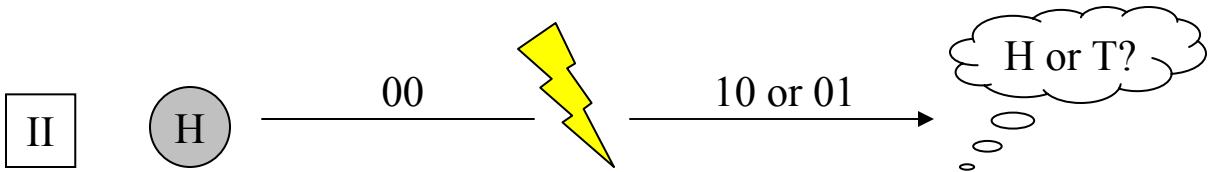
Hamming Distance

One factor to consider when transmitting a message is the possibility that the message might get corrupted in transit. For example, say we're transmitting the outcome of a coin flip over the internet. There are two possible outcomes – heads or tails.

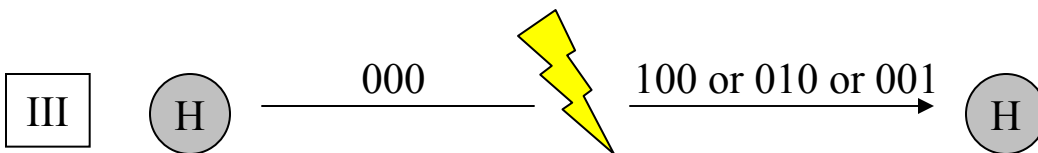
One possible encoding uses only one bit: heads = “0” and tails = “1”. We tell our friend that the coin flipped heads by transmitting a “0.” However, a routing wire is struck by lightning, changing the bit to a “1.” Our friend will incorrectly think tails.



If we increase the number of bits in the encoding to two bits (heads = 00 and tails = 11), we can detect single-bit errors. However, we still cannot correct them.



With three bits, however, we *can* detect single and two-bit errors, and we can correct single bit errors (heads = 000 and tails = 111).



For two encodings of the *same* length, the Hamming Distance (**D**) is the number of bits you need to change to turn one into the other. We increased the Hamming Distance by one in each step of this example. (I) $D = 1$ (II) $D = 2$ (III) $D = 3$.

If you have a Hamming Distance of **D**, then you can:

DETECT $(D - 1)$ bit errors

and

CORRECT $\frac{(D - 1)}{2}$ bit errors

Information

A message carries information if it tells you something new. The higher the probability of receiving that message, the less information that message carries. There are three basic equations you need to know:

Non-equivalent probability for all choices

$$\# \text{ bits of information} = \log_2 \left[\frac{1}{p_i} \right]$$

Equivalent probability for all choices

$$\# \text{ bits of information} = \log_2 \left[\frac{\# \text{ choices before}}{\# \text{ choices after}} \right]$$

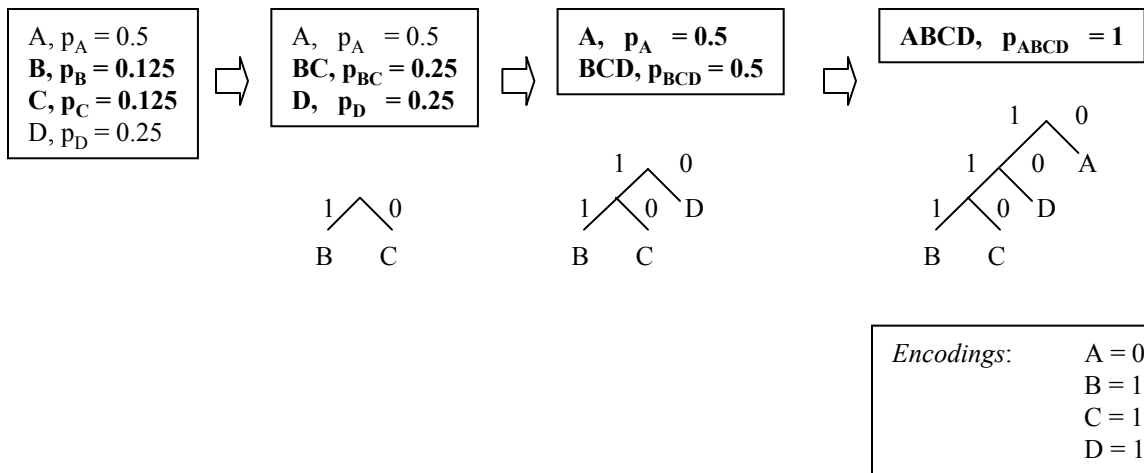
$$\text{Average \# bits of information} = \sum p_i \log_2 \left[\frac{1}{p_i} \right]$$

Huffman Encoding

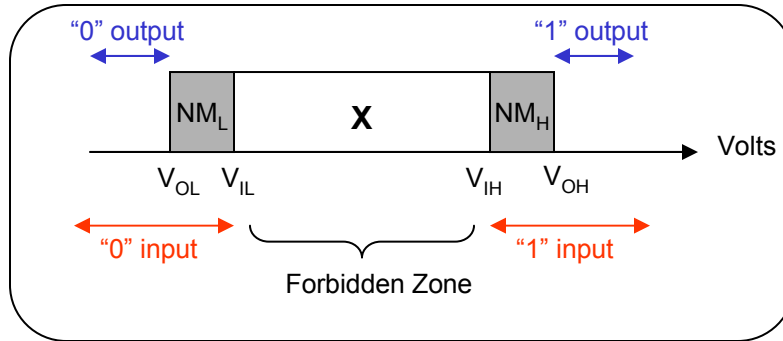
We can encode symbols using a sequence of bits. Huffman encoding is merely one of the ways we can assign a symbol to a sequence of bits. The algorithm proceeds as follows:

- (1) Choose the two members/subtrees with the lowest probability
- (2) Label one a “0” and the other a “1”
- (3) Connect them into a new subtree, and add their probabilities to produce the probability of the new subtree
- (4) Repeat 1-3 until you’re left with one big tree

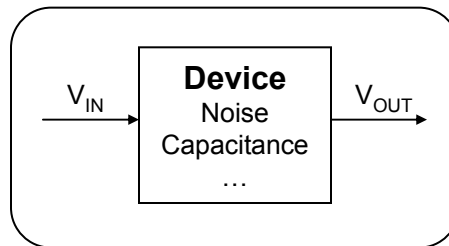
This example illustrates Huffman encoding for the symbols A, B, C, and D with their corresponding probabilities:



The Digital Abstraction



In the digital world, a voltage can represent a valid “0”, a valid “1”, or an invalid signal “X” (a voltage in the forbidden zone). Voltage is continuous, therefore we define boundaries to indicate whether it’s a 0, 1, or X.



The world is not ideal, however, so we must ensure that a “0” will never be mistaken for anything else. This is done by applying stricter boundaries to the outputs of combinational logic blocks (V_{OL} and V_{OH}) and more lenient boundaries on the inputs (V_{IL} and V_{IH}).

The difference between these two boundaries is called a **noise margin**. There are two noise margins: $NM_{L=LOW}$ and $NM_{H=HIGH}$. Think of NM_L as the amount of voltage it takes to turn a “0” into an “X”, and NM_H as the amount of voltage it takes to turn a “1” into an “X.”

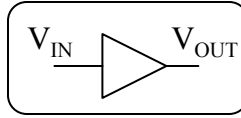
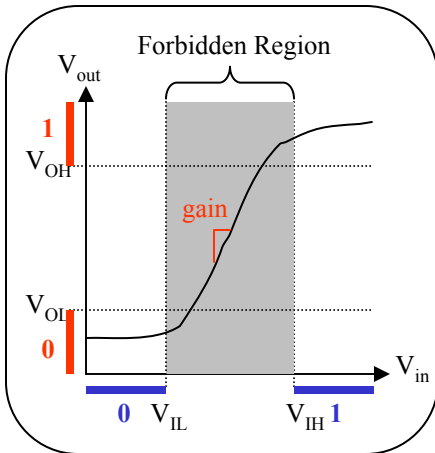
$$NM_L = V_{IL} - V_{OL}$$

$$NM_H = V_{OH} - V_{IH}$$

$$\text{Noise Margin of a combinational logic block} = NM = \min(NM_L, NM_H)$$

Combinational Devices and Static Discipline

Buffer



Static Discipline:

If $V_{IN} \leq V_{IL}$ then $V_{OUT} \leq V_{OL}$

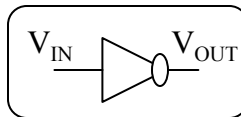
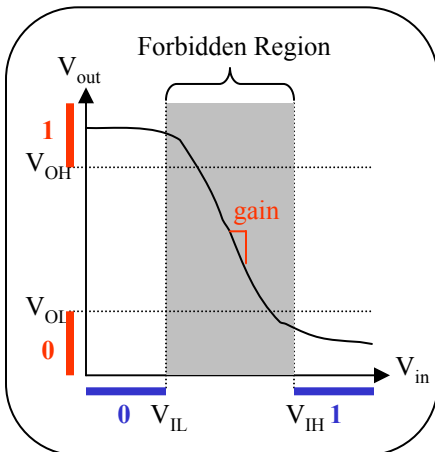
If $V_{IN} \geq V_{IH}$ then $V_{OUT} \geq V_{OH}$

A buffer basically takes in a “0” and outputs a “0”, or takes in a “1” and outputs a “1”. The buffer symbol is shown above. The graph on the left is the buffer’s Voltage Transfer Curve (VTC). The valid regions are located in the lower left and upper right corners of the curve. The shaded region in the middle is the invalid Forbidden Region.

Gain is the slope in the forbidden region of the VTC.

Note that the $|\text{gain}| > 1$ in order for the buffer to be classified as a combinational device. If it isn’t, you might get negative noise margins or the signal will degrade when you cascade the device.

Inverter



Static Discipline:

If $V_{IN} \leq V_{IL}$ then $V_{OUT} \geq V_{OH}$

If $V_{IN} \geq V_{IH}$ then $V_{OUT} \leq V_{OL}$

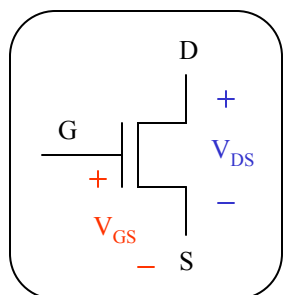
An inverter takes in a “0” and outputs a “1”, or takes in a “1” and outputs a “0”. Note that the inverter must also have $|\text{gain}| > 1$ in order to be a combinational device.

When choosing values for V_{OL} , V_{IL} , V_{IH} , V_{OH} , there is no set algorithm. The best way is to gain intuition from practice.

Some rules of thumb might include:

- (1) Pick V_{OH} and V_{OL} around the corners where the flat region meets the steep region
- (2) If you’re trying to maximize the noise margins, try to choose the highest V_{OH} and the lowest V_{OL} possible
- (3) If you’re trying to maximize the noise margins, first find a V_{OL} , V_{IL} , V_{IH} , V_{OH} that adheres to the static discipline. Then try to make one of the noise margins bigger by altering one of V_{OL} , V_{IL} , V_{IH} , V_{OH} and see if the new values adheres to the static discipline.

NFET pulldowns



A NFET can be categorized into three “analog” regions – cutoff, linear, and saturation - and two “digital” regions – on and off.

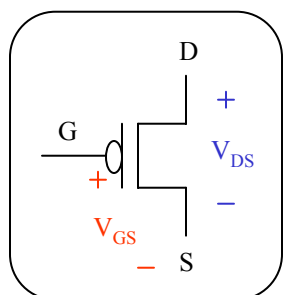
Off/Cutoff: $V_{GS} < V_{TH}$

On/Linear: $V_{GS} \geq V_{TH}, V_{DS} < V_{DSsat}$

On/Saturation: $V_{GS} \geq V_{TH}, V_{DS} \geq V_{DSsat}$

Where $V_{DSsat} = V_{GS} - V_{TH}$
 V_{TH} = a set Threshold Voltage

PFET pullups



A PFET also has three “analog” regions – cutoff, linear, and saturation - and two “digital” regions – on and off.

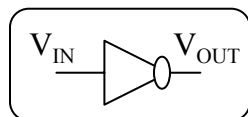
Off/Cutoff: $V_{GS} > V_{TH}$

On/Linear: $V_{GS} \leq V_{TH}, V_{DS} > V_{DSsat}$

On/Saturation: $V_{GS} \leq V_{TH}, V_{DS} \leq V_{DSsat}$

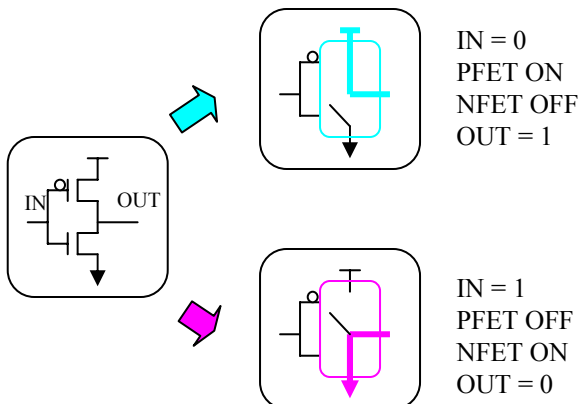
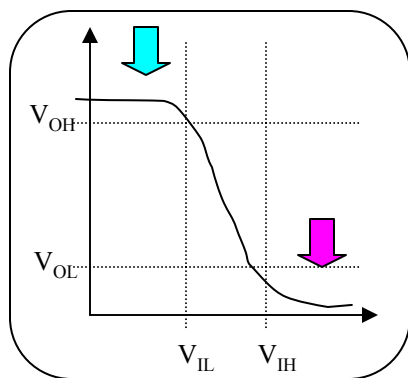
Where $V_{DSsat} = V_{GS} - V_{TH}$

CMOS Inverter

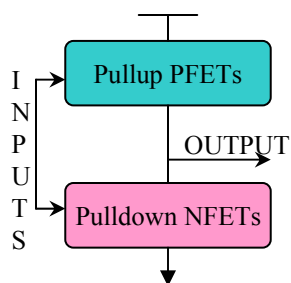


IN	OUT
0	1
1	0

A CMOS inverter has one PFET between power and the output, and one NFET between the output and ground. Think of the PFET and NFET as switches. A PFET “switch” closes when $IN=0$, and ties the output to V_{DD} . A NFET “switch” closes when $IN=1$, and ties the output to ground.



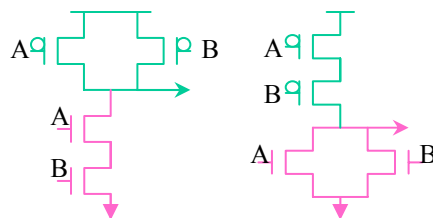
CMOS Logic



A CMOS Logic gate consists of a bunch of PFETs in the pullup (between power and the output), and a bunch of NFETs in the pulldown (between the output and ground).

The pullup may only consist of PFETs and the pulldown may only consist of NFETs. Otherwise, the logic gate will not have sufficient gain in the forbidden area. Also, since the NFETs always go in the pulldown, CMOS logic gates can only implement negative logic.

The pullup and pulldown must complement each other. If two NFETs are in parallel in the pulldown, then the corresponding PFETs must be in series in the pullup. The opposite is also true – NFETs in series must have corresponding PFETs in parallel.

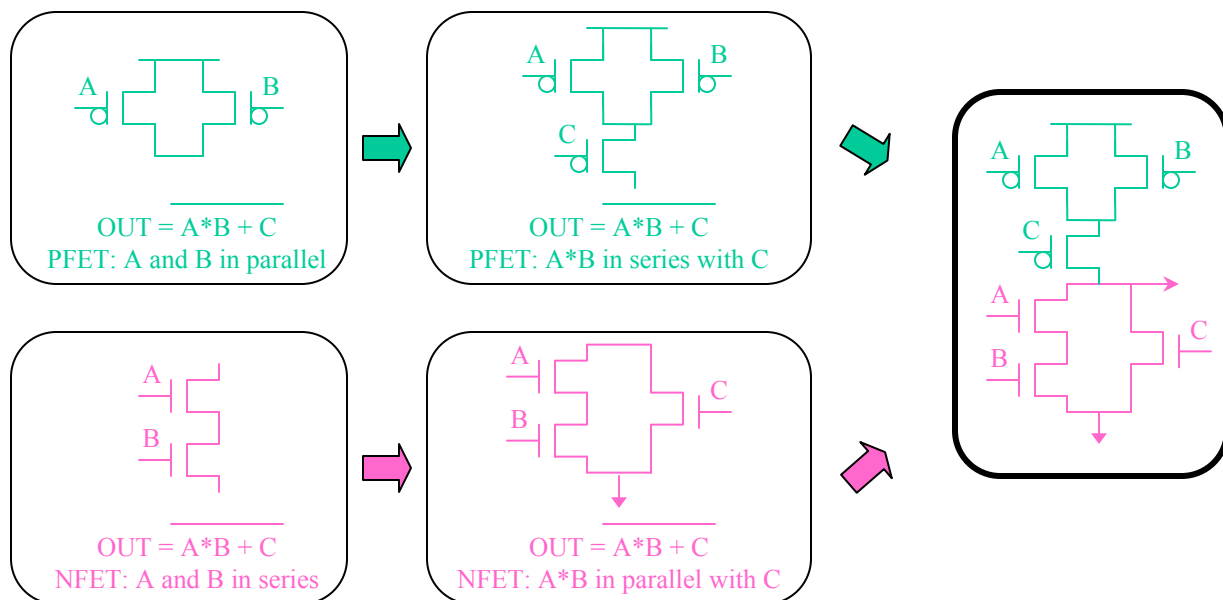


If you have a formula that you need to implement as a CMOS logic gate:

$A * B$ means A and B are in series in the pulldown and in parallel in the pullup

$A + B$ means A and B are in parallel in the pulldown and in series in the pullup

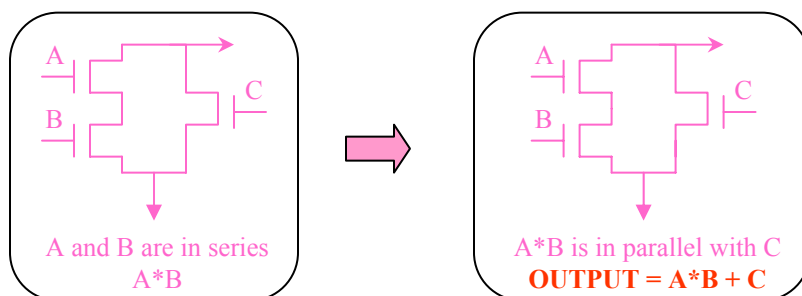
First draw the pulldown. Then after the pulldown is done, draw the pullup as the complement of the pulldown. In the example below, we must draw the CMOS logic gate for **OUTPUT = $A * B + C$** .



If you have a CMOS logic gate and you wish to derive the formula, follow these steps:

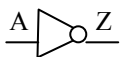
- 1) Look at the pulldown only
- 2) A and B in series means there is an “ $A*B$ ” in the equation
- 3) A and B in parallel means there is an “ $A+B$ ” in the equation
- 4) Negate the entire thing after you’re done with all the pulldown NFETs.

Step 4 is necessary because you’re dealing with NFETs. You will probably want to do a few similar practice problems to gain intuition.



CMOS Logic Gates

INVERTER



A	Z
0	1
1	0

NAND



A	B	Z
0	0	1
0	1	1
1	0	1
1	1	0

NOR



A	B	Z
0	0	1
0	1	0
1	0	0
1	1	0

XNOR



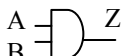
A	B	Z
0	0	1
0	1	0
1	0	0
1	1	1

BUFFER



A	Z
0	0
1	1

AND



A	B	Z
0	0	0
0	1	0
1	0	0
1	1	1

OR



A	B	Z
0	0	0
0	1	1
1	0	1
1	1	1

XOR



A	B	Z
0	0	0
0	1	1
1	0	1
1	1	0

The figures above contain eight typical standard CMOS logic gates, their symbolic representation, and their truth tables. You can combine logic gates to implement *almost* any function. However, keep in mind that negative logic (inverter, nand, nor, xnor) is usually smaller and faster than positive logic (buffer, and, or, xor).

Boolean Algebra

$A * B$ means “A and B”

$\overline{A * B}$ means “A nand B”

$A + B$ means “A or B”

$\overline{A + B}$ means “A nor B”

$A \oplus B$ means “A xor B”

$A \oplus B$ means “A xnor B”

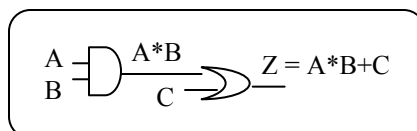
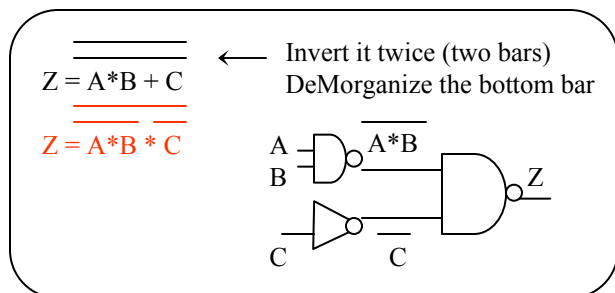
DeMorgan's Law

$$\overline{A * B} = \overline{A} + \overline{B}$$

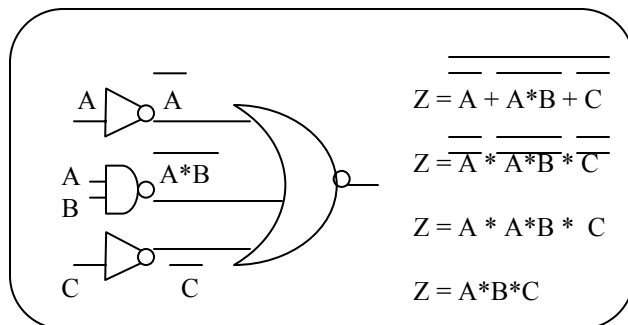
$$\overline{A + B} = \overline{A} * \overline{B}$$

Buildings Circuits with Logic Gates

If you wish to implement a function $Z = A*B + C$ using only standard logic gates such as those pictured above, you can implement it in multiple ways. The two possibilities below (all negative logic on the left and all positive logic on the right) are not necessarily faster or smaller than one another – it depends on the size and speed of the standard logic gates.



If you wish to derive a logic equation from an existing circuit, try to derive intermediate equations, starting at the inputs and working incrementally towards the output. You can use boolean algebra to manipulate the formula to whatever format you prefer.

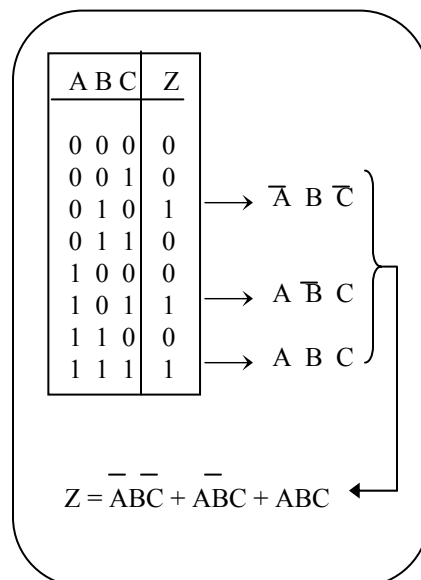


Sum of Products

Sum of Products is an easy way to derive a formula from a truth table. The formula is a sum of products. Each product represents one row where the output $Z = 1$. For each product/row, write down all the inputs and put individual bars over the inputs that equal zero.

If there are fewer $Z = 0$ than $Z = 1$ outputs, you can find the sum of products using the same method, but use the $Z = 0$ rows as the products and negate the entire sum of products.

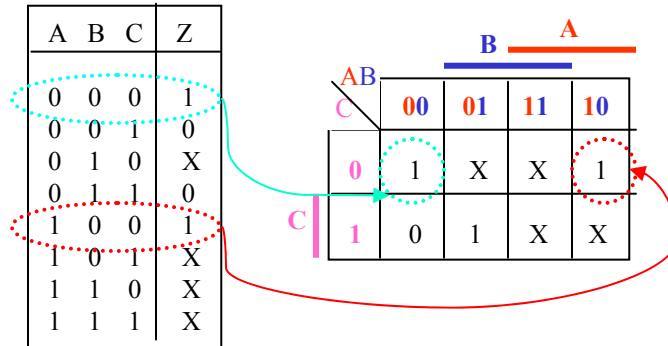
From here, you should be able to find the *minimal* sum of products by reducing the formula using boolean arithmetic. This can be quite tedious. An easier method for finding the minimal sum of products is presented in a subsequent section on Karnaugh Maps.



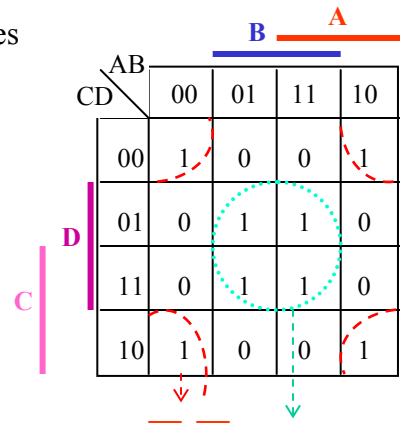
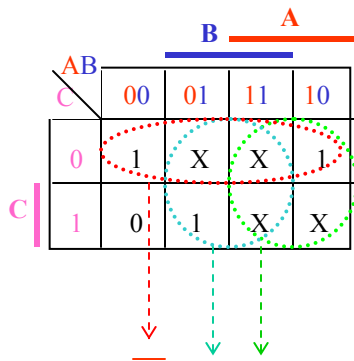
Karnaugh Maps

Karnaugh Maps are extremely useful for finding the minimal sum of products.

- 1 Translate a truth table to a Karnaugh Map.
→ Note the order of the AB field – neighbors must differ by one bit



- 2 Find and circle all 1's
 - Don't Cares 'X' can be circled, if needed
 - Circles must have dimensions that are a power of 2 (1, 2, 4, 8, etc.)
 - You can overlap circles
 - The Karnaugh Map folds over, so circles can go off the sides
 - Look for the largest groups of 1's
 - Try to get the smallest number of circles

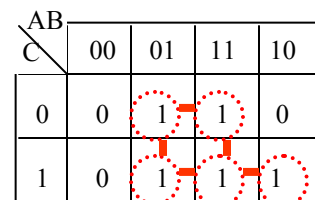


- 3 $OUT = C + B + A$

$$OUT = B D + B D$$

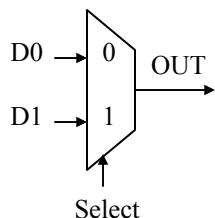
Lenient Combinational Logic

Take a karnaugh map that is directly derived from the circuit. Each gap between neighboring circles is an opportunity for a glitch. The circuit is lenient if there are no gaps.



Mux

A d-input mux has a selector input, d data inputs (d is a power of two), and one output.

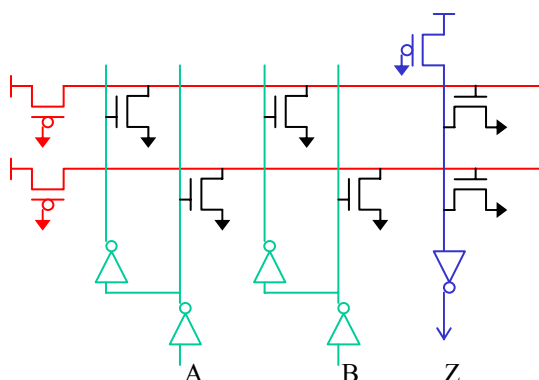


D0	D1	Select	OUT
0	X	0	0
1	X	0	1
X	0	1	0
X	1	1	1

If Select = 0, then OUT = D0

If Select = 1, then OUT = D1

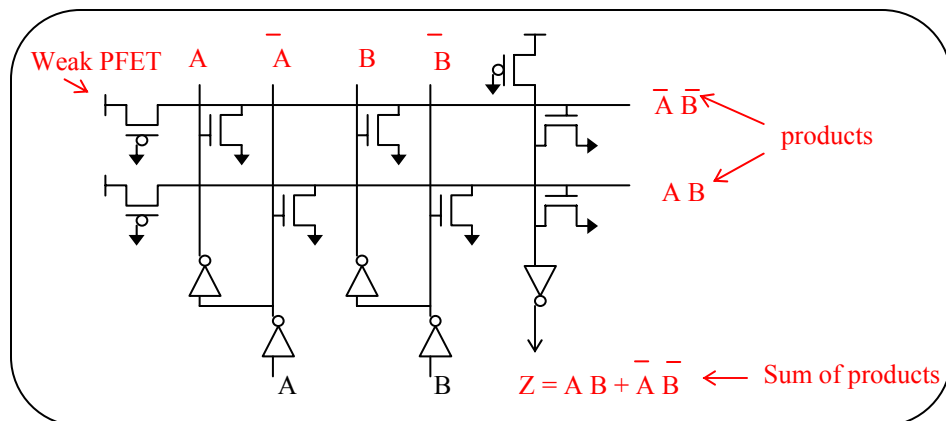
ROM / PLA



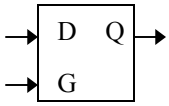
A Read-Only Memory and a Programmable Logic Array consist of **address input lines**, **word lines**, and **bit lines**. You can add NFETs to configure the ROM or PLA to perform many functions. Each word line and bit line is driven by a weak PFET whose input is tied to ground. The weak PFET serves to drive the wire high when there is no path to ground, but is easily overpowered by a NFET.

A ROM specifies all combinations of inputs. A PLA looks like a small ROM because unnecessary FETs and wires are eliminated. If you are given the PLA pictured above, follow these steps to figure out the formula for the output:

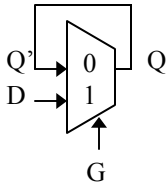
1. Write the logic equivalent of each word line at the top
2. Each bit line driving the gate of a NFET is a *product* in the sum of products
3. For each NFET drain attached to the bit line, negate the gate input.
Repeat to find product.
4. $Z = \text{OR of all the products}$



Latch



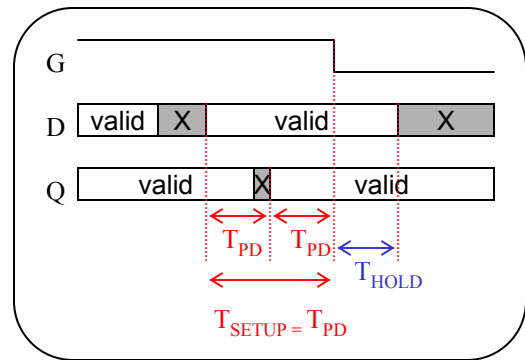
A latch has one enable input G , one data input D , and one data output Q . Think of a latch as a gate. If the latch is enabled, the gate opens and the input data propagates through to the output. If the latch is not enabled, then the gate is closed, and the output remains the same.



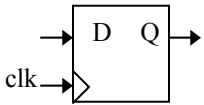
One way to implement the latch is to use a two-input mux. This introduces some timing constraints on the inputs of the latch. The data input D must be valid T_{SETUP} before G falls and D must stay valid for T_{HOLD} after G falls.

$T_{\text{SETUP}} = 2 T_{\text{PD}}$ because ...

- (1) One T_{PD} for D to propagate through to Q
- (2) Hold Q' one T_{PD} for mux to promise a stable output after G falls

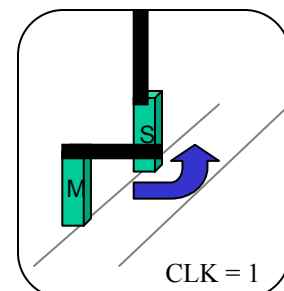
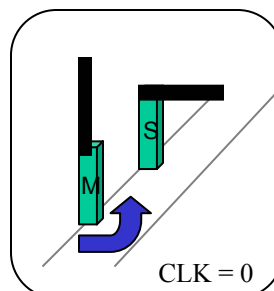
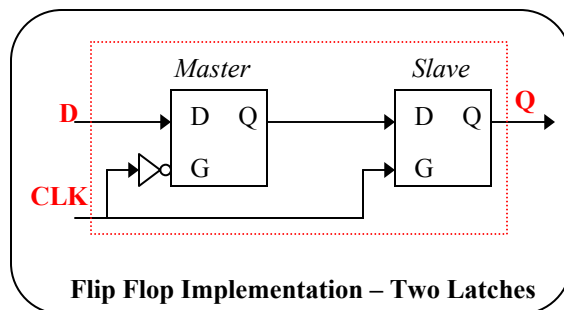


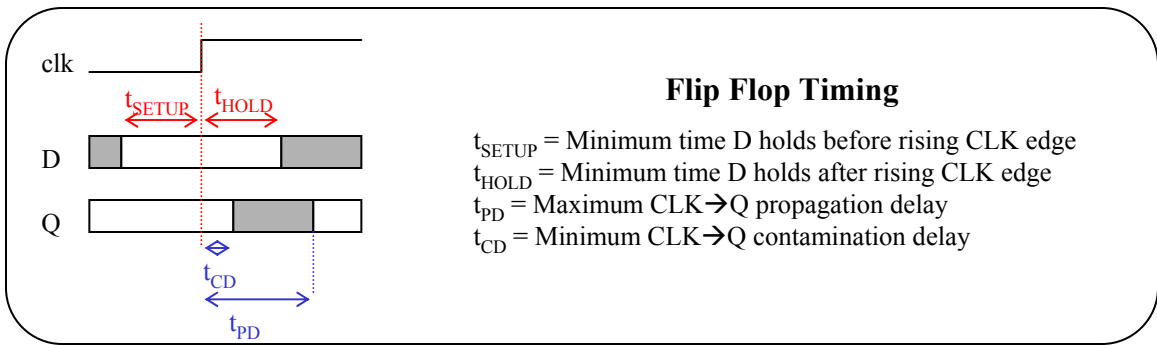
Flip Flop



A flip flop, which is interchangeably referred to as a register, has a clock input, a data input D , and a data output Q . D propagates to Q *only at the rising 0 \rightarrow 1 edge of the clock*. Otherwise, the output stays the same.

A flip flop is implemented with two latches – a master and a slave. As in the 6.004 lecture, think about this implementation as two tollbooth gates that open one at a time. The master passes the data into the middle when the clock is low. When the clock goes high, the master immediately closes to trap the data, and the slave opens to let the trapped data through to the output.

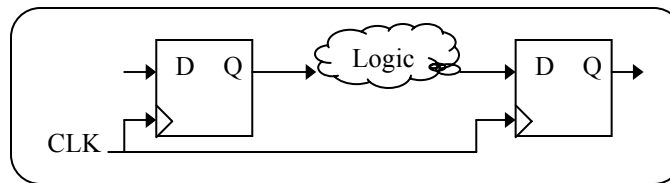




$t_{\text{CD,MASTER}} > t_{\text{HOLD,SLAVE}}$ to ensure that data will not pass through both the master and the slave while both are transitioning. If this inequality does not hold, then you can ensure that the flip flop works by inserting a buffer between the master and slave such that

$$t_{\text{CD,MASTER}} + t_{\text{CD,BUFFER}} > t_{\text{HOLD,SLAVE}}$$

Flip Flop Circuits



The circuit contains register 1 on the left, some random combinational logic, and register 2 on the right. There are two inequalities that must hold in order for this circuit to work.

$$t_{\text{CD,REG1}} + t_{\text{CD,LOGIC}} > t_{\text{HOLD,REG2}}$$

to ensure that data coming out of reg1 one does not corrupt the data going in to reg2.

$$t_{\text{CLK}} \geq t_{\text{PD,REG1}} + t_{\text{PD,LOGIC}} + t_{\text{SETUP,REG2}} \quad \text{where } t_{\text{CLK}} = \text{the clock period}$$

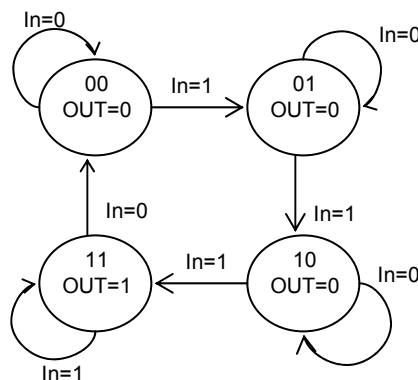
to ensure that reg2's input will meet reg2's setup time.

Finite State Machine

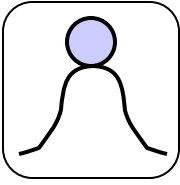
Current state			Next state		Output Based on Current State
s1	s0	IN	s1'	s0'	OUT
0	0	0	0	0	0
0	0	1	0	1	0
0	1	0	0	1	0
0	1	1	1	0	0
1	0	0	1	0	0
1	0	1	1	1	0
1	1	0	0	0	1
1	1	1	1	1	1

An FSM must designate a transition for every possible input. Also, the arcs leaving a state must not allow two choices for the same input.

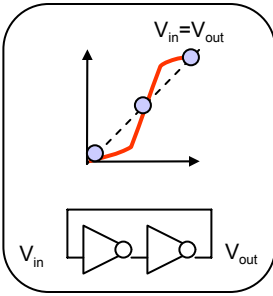
A finite state machine (FSM) contains input signals, state signals, and output signals. The example FSM pictured here contains one input, two state signals s1 and s0, and one output. The circles contain the states for s1 and s0 and tell us whether we're in state 00, 01, 10, or 11. The arcs tell us the next state based on the input.



Metastability



Imagine a ball at the top of a hill. If the hill is not flat, there are three possible states – a stable state at the bottom left, a stable state at the bottom right, and a metastable state at the top of the hill. That ball is in a metastable state because (1) the top of the hill is an unstable state – a small perturbation will make it roll down the hill, and (2) the ball will *eventually* settle to a stable state at the bottom of the hill.

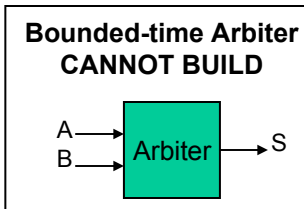


A metastable state in the digital world has the following properties:

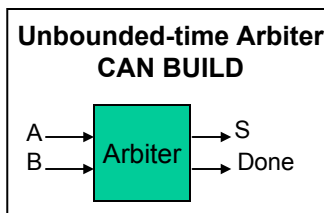
- 1) Invalid
- 2) Unstable – a small perturbation changes it to a 0 or 1
- 3) It will settle in *unbounded* time
- 4) The longer you wait, the more likely it has stabilized

The two-inverter configuration has three possible states where $V_{in} = V_{out}$: two stable states at valid 1 and 0, and one metastable state in the forbidden region.

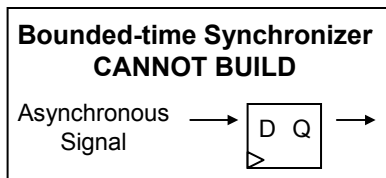
You cannot build certain devices with perfect reliability because they *cause* metastability problems. Some example systems are shown below:



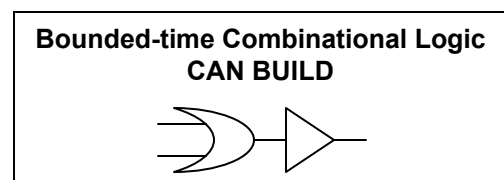
Sample *S* at a *specified time* after *A* and *B* have risen to find out which rising edge came first – *A* or *B*. The Arbiter may enter a metastable state if *A* and *B* rise at exactly the same time. The Arbiter will eventually come to decision in *unbounded* time but might be metastable and invalid when *S* is sampled.



Sample *S* when the *Done* signal is high to find out which rising edge came first – *A* or *B*. *S* may enter a metastable state if *A* and *B* rise at exactly the same time. However, we are guaranteed that *S* will be valid when *Done* is high, therefore *this* arbiter will not cause other metastable states.

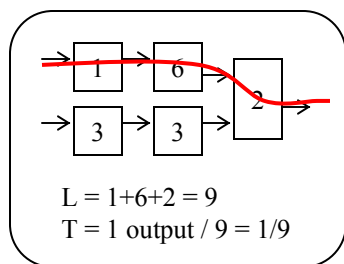


If the *D* input does not meet the flip flop's setup and hold time, the *Q* output might be metastable.



We are guaranteed a valid output after a specified propagation delay.

Unpipelined Circuits

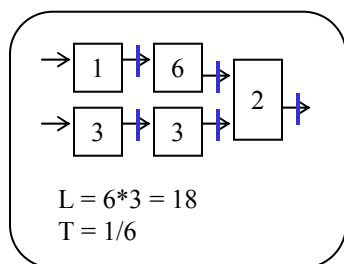


Latency = L = time an input takes to appear at the output

Throughput = $T = \# \text{outputs} / \text{time}$

The circuit contains multiple blocks of combinational logic, each with a $t_{CD} = 0$ and a t_{PD} indicated by the number inside the box. The latency L of the circuit is the longest path from the input to the output. The entire circuit is combinational logic, so the worst-case throughput T is 1 output / L .

Pipelined Circuits



Pipelined circuits contain **registers** in-between the logic blocks. For simplicity, we assume that these are ideal registers with zero contamination and propagation delay, and zero setup and hold time. The register's clock cycle t_{CLK} is equal to the propagation delay of the slowest block.

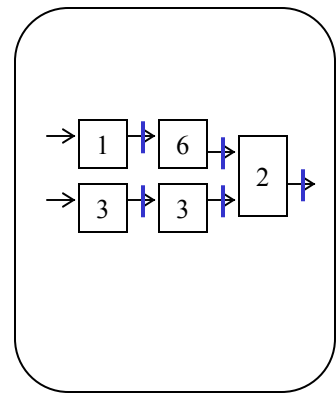
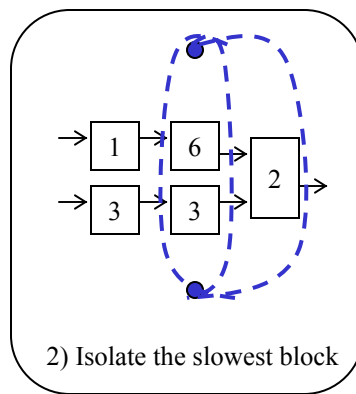
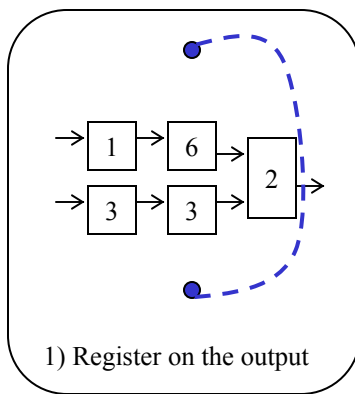
Latency = $L = K * t_{PD, \text{slowest stage}}$ where $K = \# \text{stages} = \# \text{registers on each path}$

Throughput = $T = \frac{1}{t_{PD, \text{slowest stage}}}$

Pipelining a circuit typically adds latency and increases throughput. In the example above, the 6 ns clock cycle forces an input to spend 6 ns in every stage. Since there are three stages, the input spends 18 ns before its corresponding result appears at the output as opposed to the unpipelined circuit's 9 ns latency. However, pipelining the circuit ensures that a new result appears at the output every 6ns, thus increasing the throughput to 1/6.

Use the following rules of thumb to pipeline a circuit:

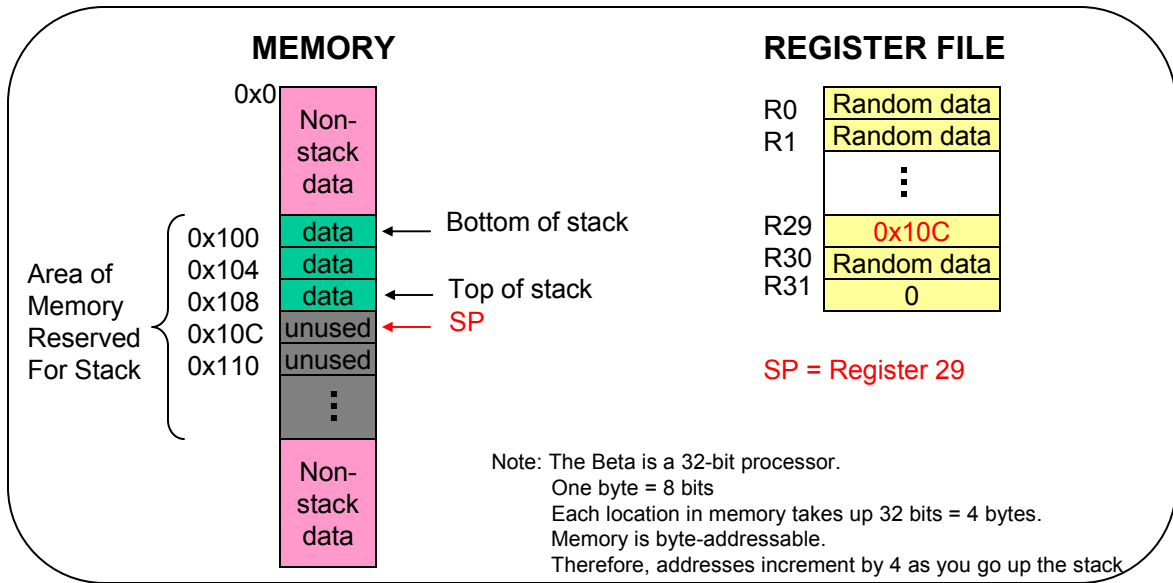
- 1) Put two dots on opposite sides of the circuit. Draw lines from one dot to another
- 2) Always put a register on the output
- 3) Isolate the slowest block with registers. Make this block the slowest stage for maximum throughput
- 4) Never cross a path twice with the same line. This helps ensure that each path has the same number of registers
- 5) Add a register every time a line crosses a wire.



The example above is one of many ways to pipeline this circuit. Any configuration is valid as long as there are an equal number of registers on each path from input to output.

Intro to Stacks

A stack is basically a stack of data located in memory. The bottom of the stack is located at a set address. The top of the stack is designated by the stack pointer (SP, Register 29), a register that contains the address of the next *unused* memory location.

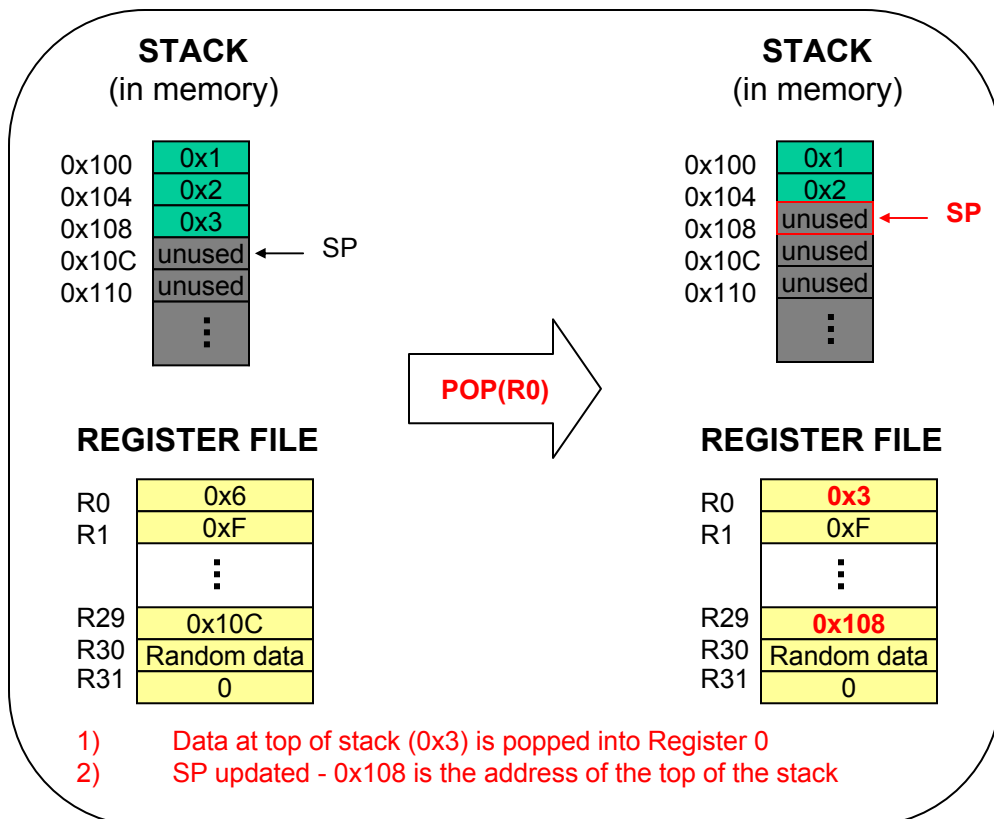
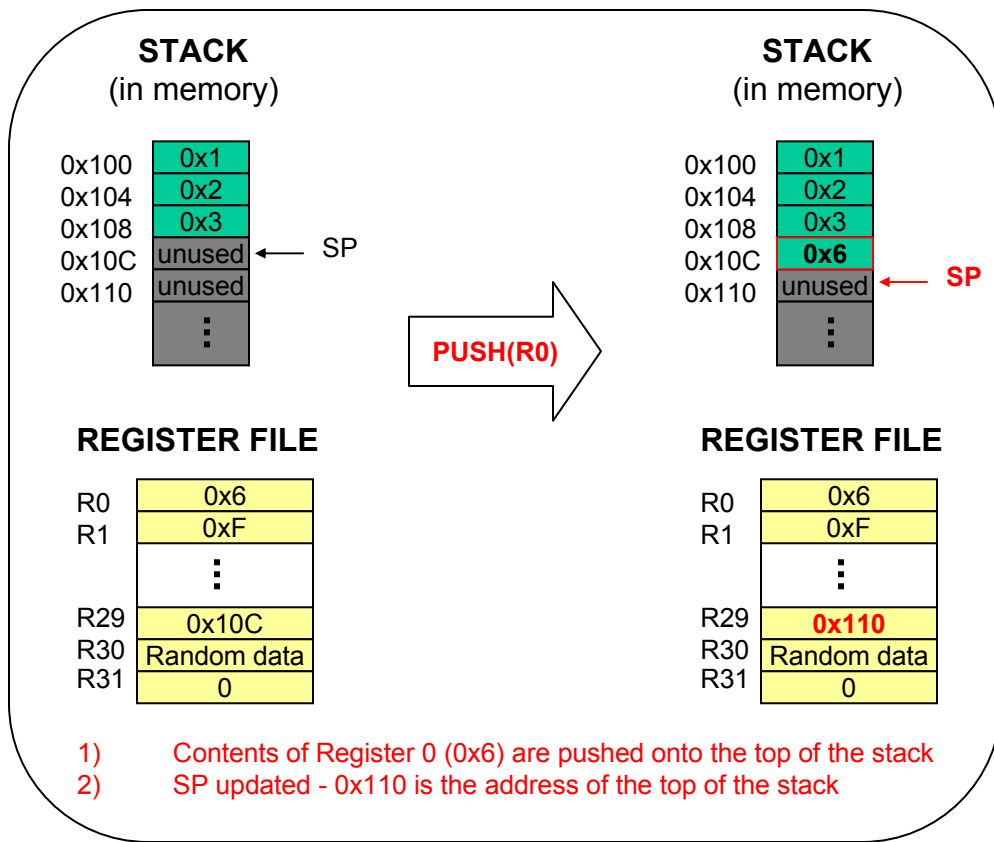


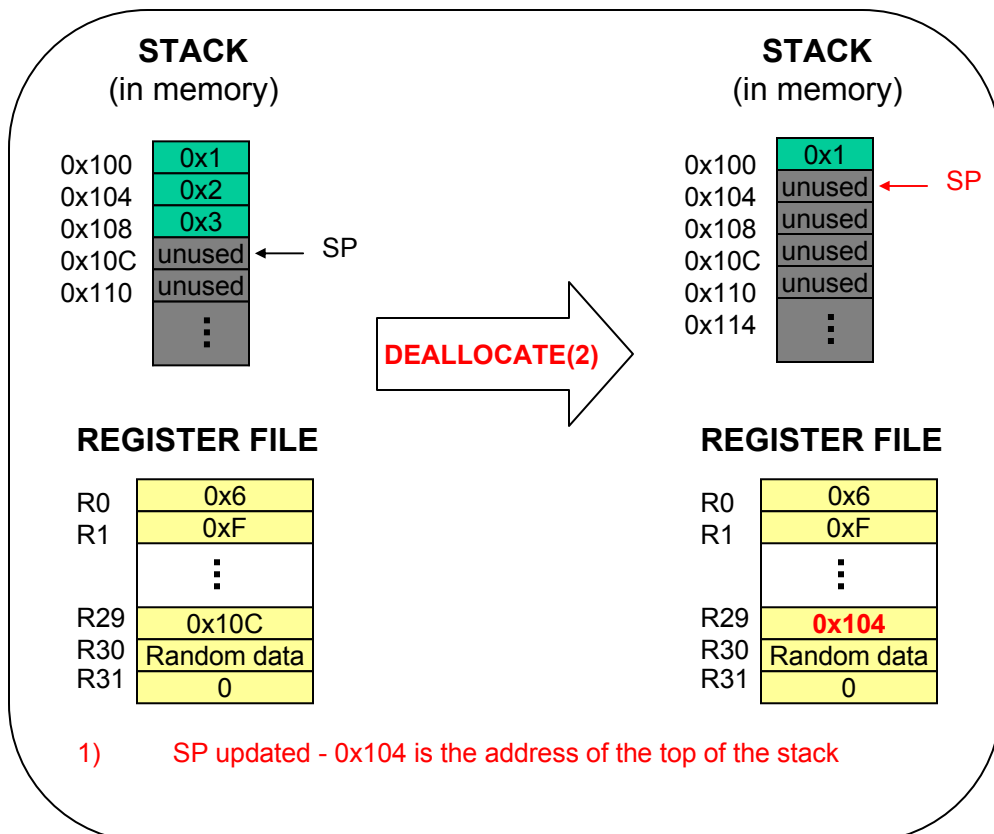
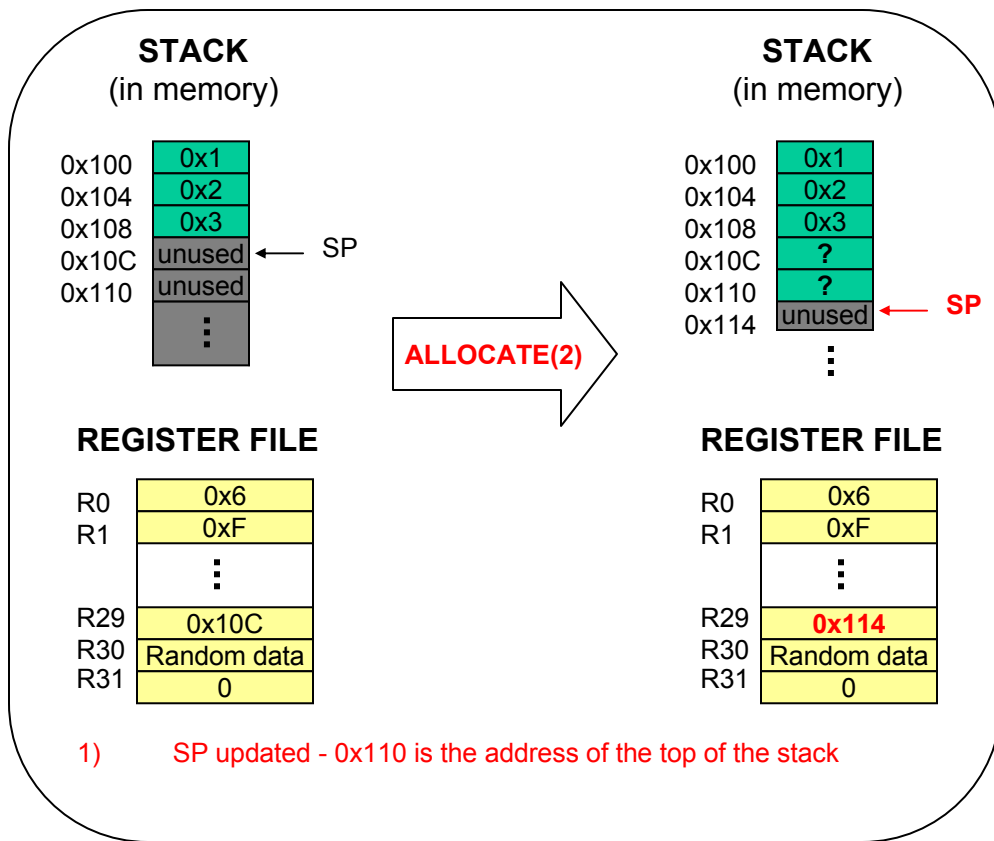
Stack Macros

There are four stack macros you should be familiar with:

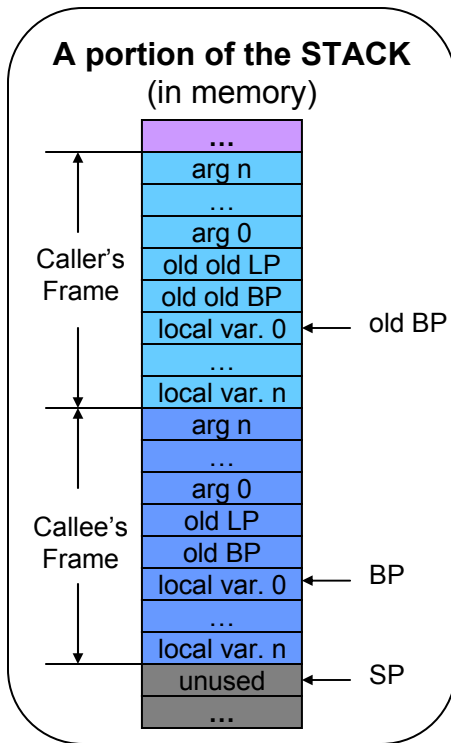
- | | |
|----------------------|--|
| PUSH(Rx) | (1) Push contents of register Rx on top of the stack
(2) $SP \leftarrow SP + 4$ |
| POP(Rx) | (1) Pop contents of register Rx on top of the stack
(2) $SP \leftarrow SP + 4$ |
| ALLOCATE(k) | (1) Reserve k lines: $SP \leftarrow SP + 4*k$ |
| DEALLOCATE(k) | (1) Release k lines: $SP \leftarrow SP - 4*k$ |

Note that PUSH and POP are each *two instruction* macros – one to push/pop between the register and the stack, and one to update the Stack Pointer (SP).





Stack Frames



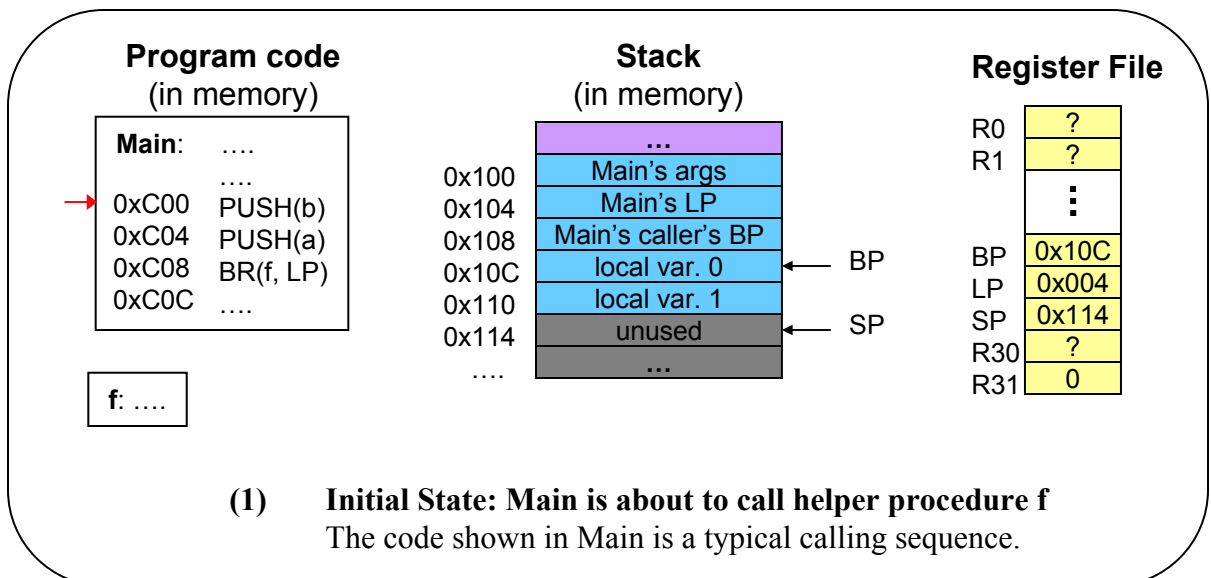
Stack frames allow procedures to call other procedures (**caller** calls the **callee**), and prevent procedures from corrupting each other's data.

Two new registers are added to support Stack Frames. Register 27, Base Pointer **BP**, points to local variable 0 in the frame and is used to access the procedure's arguments and local variables. Register 28, Link Pointer **LP**, contains the return address of the caller. As stated in previous sections, Register 29 is the Stack Pointer **SP**.

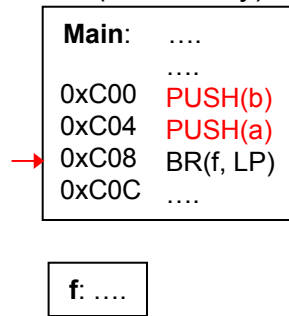
A stack frame consists of (1) arguments pushed in reverse order, (2) the caller's link pointer, (3) the caller's base pointer, and (4) a bunch of local variables.

Two-level stack example

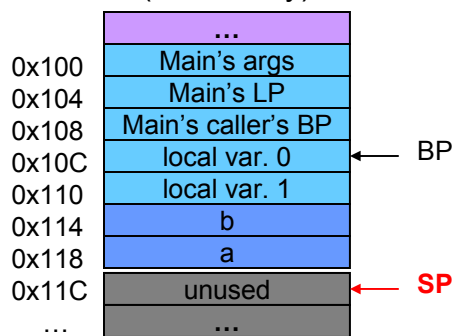
In this example, we have a caller **main** and a callee **f**. The initial state of the stack is pictured below. Main has been running for awhile and is about to call **f**. The program code arrow points to the line that is about to be executed. To conserve space, only relevant portions of the program code are shown in each snapshot.



Program code (in memory)



Stack (in memory)

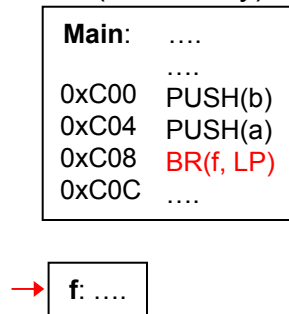


Register File

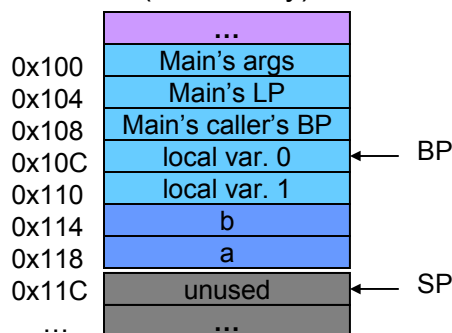
R0	?
R1	?
	⋮
BP	0x10C
LP	0x004
SP	0x11C
R30	?
R31	0

(2) Main pushed f's arguments onto the stack in reverse order

Program code (in memory)



Stack (in memory)

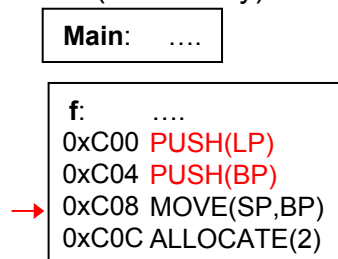


Register File

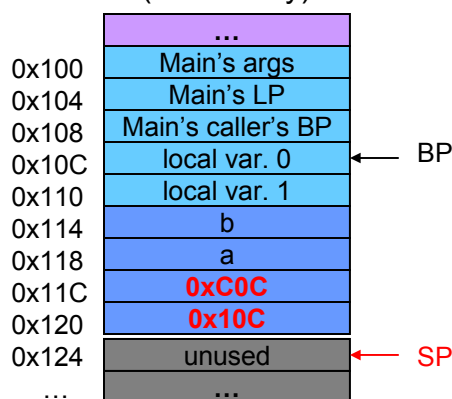
R0	?
R1	?
	⋮
BP	0x10C
LP	0xC0C
SP	0x11C
R30	?
R31	0

(3) Main branched to f and stores return address in LP

Program code (in memory)



Stack (in memory)



Register File

R0	?
R1	?
	⋮
BP	0x10C
LP	0xC0C
SP	0x124
R30	?
R31	0

(4) f pushed the return address (LP) and Main's BP
 The code shown in f is a typical entry sequence.

Program code (in memory)

Main:

f:
0xC00 PUSH(LP)
0xC04 PUSH(BP)
0xC08 **MOVE(SP,BP)**
→ 0xC0C **ALLOCATE(2)**

Stack (in memory)

	...
0x100	Main's args
0x104	Main's LP
0x108	Main's caller's BP
0x10C	local var. 0
0x110	local var. 1
0x114	b
0x118	a
0x11C	0xC0C
0x120	0x10C
0x124	unused
...	...

BP

SP

Register File

R0	?
R1	?
	⋮
BP	0x124
LP	0xC0C
SP	0x124
R30	?
R31	0

(5) f updated BP to reflect the current stack frame

Program code (in memory)

Main:

f:
0xC00 PUSH(LP)
0xC04 PUSH(BP)
0xC08 MOVE(SP,BP)
0xC0C **ALLOCATE(2)**
→ 0xC10 ...

Stack (in memory)

	...
0x100	Main's args
0x104	Main's LP
0x108	Main's caller's BP
0x10C	local var. 0
0x110	local var. 1
0x114	b
0x118	a
0x11C	0xC0C
0x120	0x10C
0x124	reserved
0x128	reserved
0x12C	...

BP

SP

Register File

R0	?
R1	?
	⋮
BP	0x124
LP	0xC0C
SP	0x12C
R30	?
R31	0

(6) f allocated space for two local variables

Program code (in memory)

Main:

f:

→ 0xC30 MOVE(val, R0)
0xC34 MOVE(BP, SP)
0xC38 POP(BP)
0xC3C POP(LP)
0xC40 JMP(LP)

Stack (in memory)

	...
0x100	Main's args
0x104	Main's LP
0x108	Main's caller's BP
0x10C	local var. 0
0x110	local var. 1
0x114	b
0x118	a
0x11C	0xC0C
0x120	0x10C
0x124	Local var 0 ← BP
0x128	Local var 1
0x12C	... ← SP

Register File

R0	Z
R1	?
	⋮
BP	0x124
LP	?
SP	0x12C
R30	?
R31	0

(7) f did some computing and put the return value Z in R0

f loaded argument a into a register using: LD(BP, -12, Reg)

f loaded argument b into a register using: LD(BP, -16, Reg)

f loaded local variable 0 into a register using: LD(BP, 0, Reg)

f loaded local variable 1 into a register using: LD(BP, 4, Reg)

Program code (in memory)

Main:

→ 0xC08 BR(f, LP)
0xC0C

f:

0xC30 **MOVE(val, R0)**
0xC34 **MOVE(BP, SP)**
0xC38 **POP(BP)**
0xC3C **POP(LP)**
0xC40 **JMP(LP)**

Stack (in memory)

	...
0x100	Main's args
0x104	Main's LP
0x108	Main's caller's BP
0x10C	local var. 0 ← BP
0x110	local var. 1
0x114	b
0x118	a
0x11C	... ← SP

Register File

R0	Z
R1	?
	⋮
BP	0x10C
LP	0xC0C
SP	0x12C
R30	?
R31	0

(8) f executed a return sequence shown above:

0xC30 - Stored return value into memory

0xC34 - Deallocated local variables by moving SP to BP

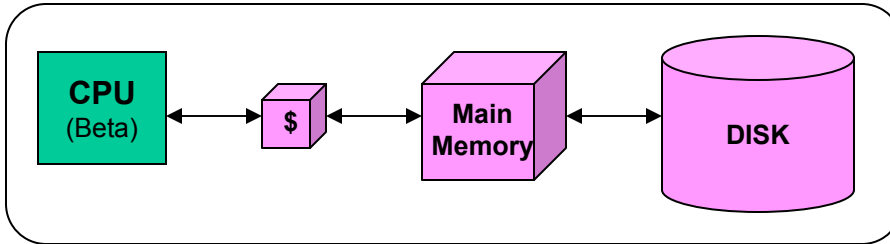
0xC38 - Restored the BP

0xC3C - Restored the LP

0xC40 - Jumps to return address stored in LP

Memory Hierarchy

Memory is one of the bottlenecks of a pipelined processor. We can alleviate this bottleneck with a hierarchy – a small fast cache (often notated as “\$”), a larger slower main memory, and an enormous but *very* slow disk. Main memory has everything the cache has and more. Similarly, disk has everything contained in main memory and more.



When the CPU accesses a memory location A, it goes through these steps:

- (1) *Check Cache* – If A is in the cache (cache hit), then return $DATA_A$, else go to step 2
- (2) *Check Main Memory* – If A is in memory, then return $DATA_A$ and put it in the cache, else go to step 3.
- (3) A will always be in disk. Put A in main memory and the cache, and return the value.

For simplicity, we will assume (for now) that we never have to go to disk – if there is a cache miss on A, we will always find A in main memory.

Caches are small, fast memory that contain temporary copies of selected memory locations. For 6.004, you can assume that the CPU runs at a clock speed where a cache hit on a memory access will return the cache data in one cycle, and a cache miss will cause the CPU to stall until the data is ready from memory. If the miss rate isn't too high, using a cache typically improve the average memory access time.

$$t_{\text{ave}} = t_{\text{cache}} + (\text{miss rate}) * t_{\text{memory}}$$

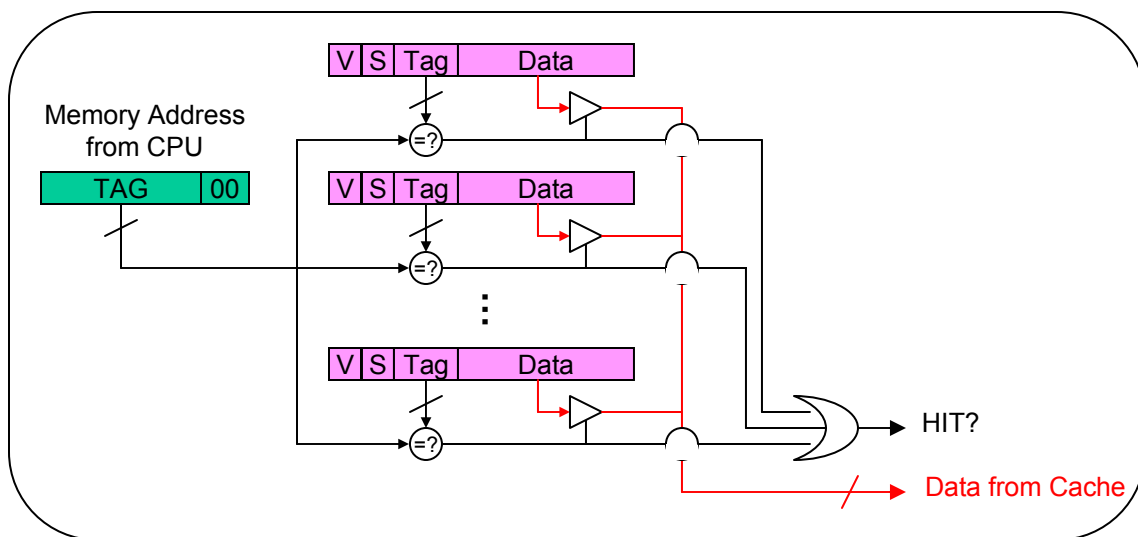
t_{cache} = time it takes to access the cache

t_{memory} = time it takes to access main memory

Temporal Locality

Caches are useful because they exploit temporal locality in memory accesses. If you access memory location A, you're more likely to access A again in the near future. Every time there is a cache miss, that memory location will go in the cache and replace something that hasn't been used for awhile. Spatial locality will be discussed in a later section.

Fully Associative Cache



A cache contains many cache lines. A cache line contains (1) a valid bit indicating if the data in that line is valid, (2) optional status bits that indicate if the data is dirty, read-only, etc., (3) a tag formed from the memory address, and (4) the data corresponding to that address.

A fully-associative cache is one of the three cache architectures described in this handbook. It differs from other cache architectures in that it *can put data anywhere*.

When a memory address is presented to a fully associative cache, the address tag is compared to tags from all the valid lines in the cache. If there is a match, it's a cache hit and the data from that line is sent to the CPU. Otherwise it's a cache miss.

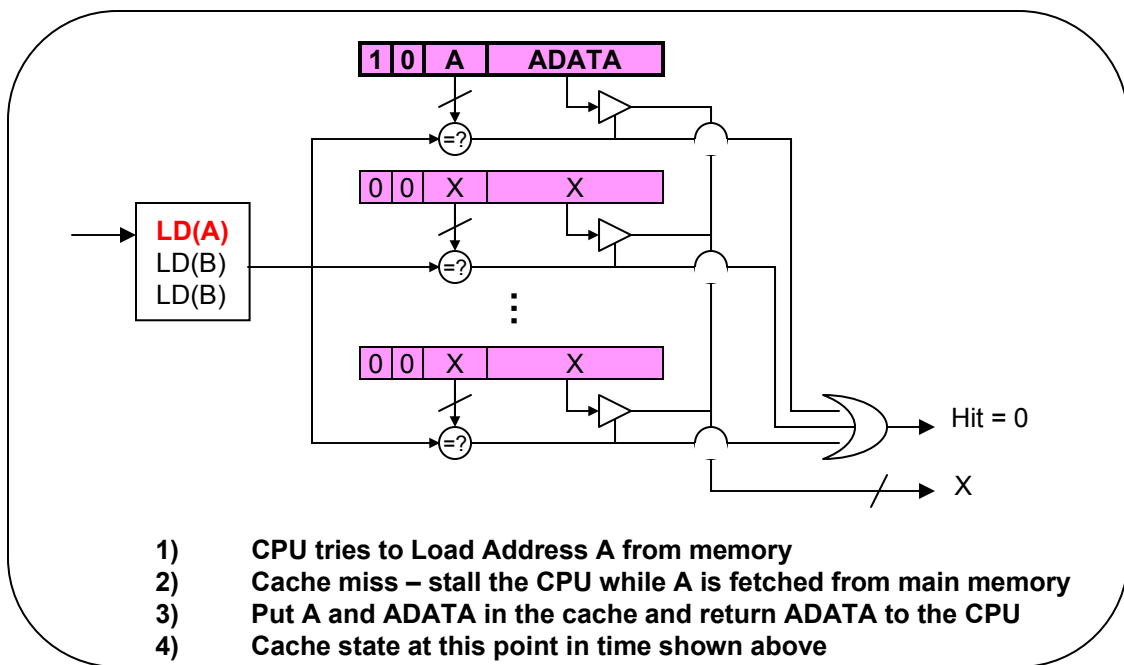
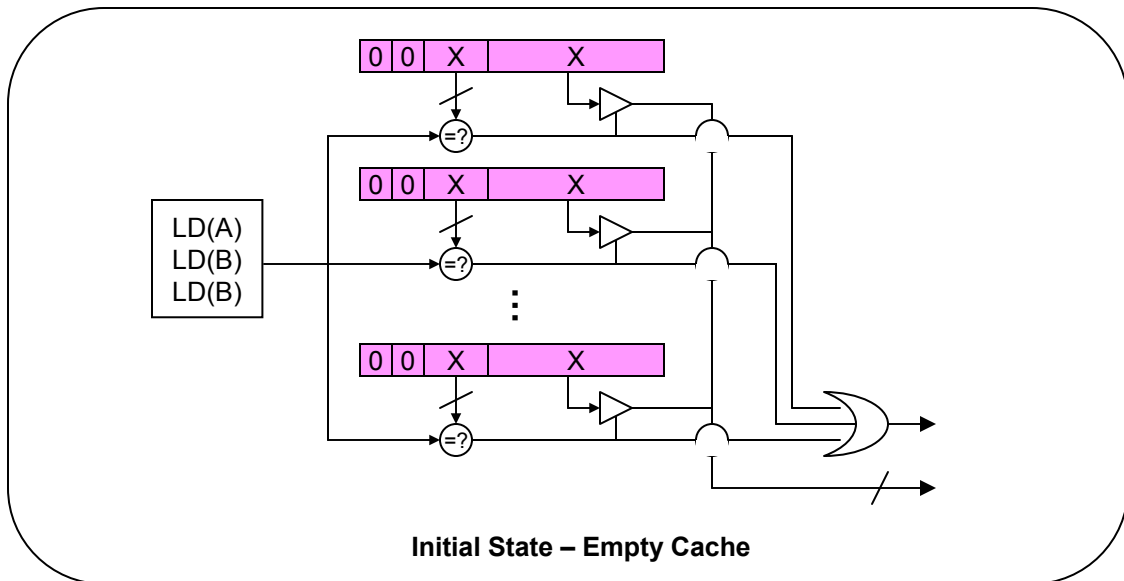
If there is a cache miss when the cache is full (all the lines are valid), we figure out which line to replace using one of the replacement policies outlined in a later section. If the line that is replaced contains valid and dirty data ($V=1$ and $D=1$), then we also have to write the data back to main memory. If it's not dirty, we just throw the line out.

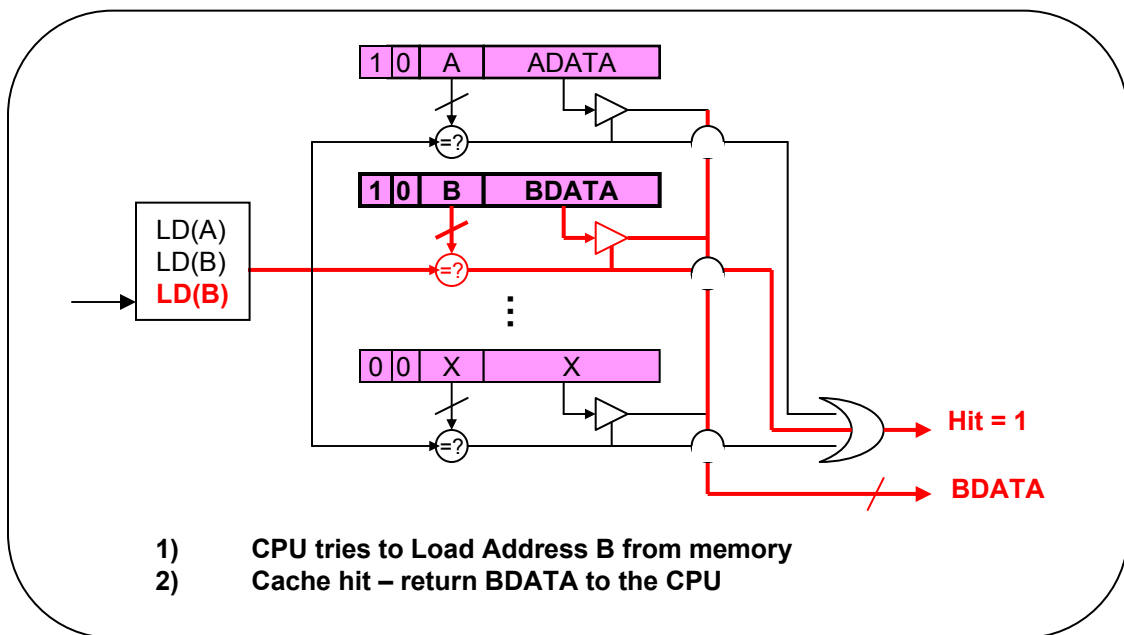
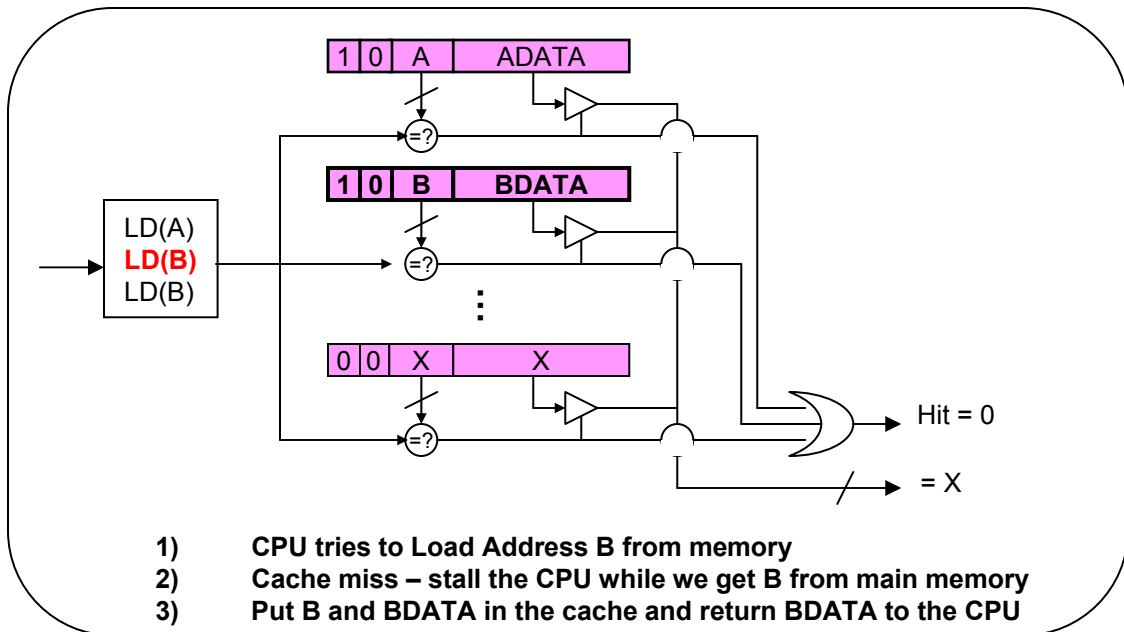
Some useful equations regarding fully associative caches:

$$\text{Capacity} = \# \text{ lines} * (1 \text{ valid bit} + S \text{ status bits} + T \text{ tag bits} + D \text{ data bits})$$

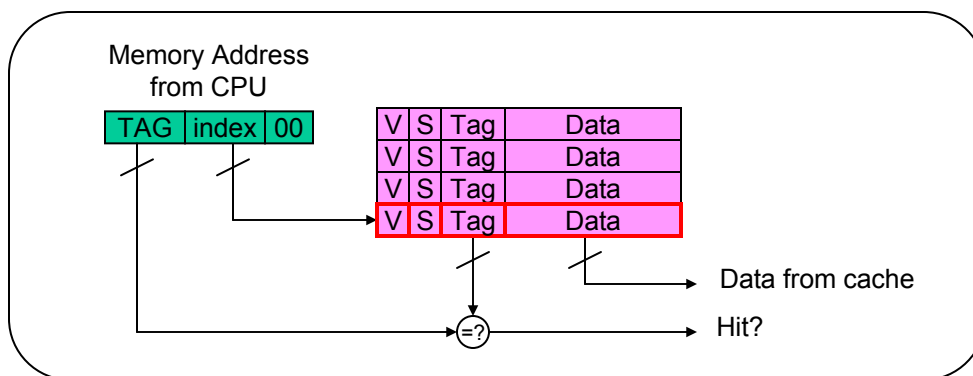
Typical values for the 6.004 BETA: $T = 30$ bits, $D = 32$ bits

The following example illustrates three memory accesses, starting with an empty fully associative cache. Each picture illustrates the state of the cache *after* the things in the commentary have executed.





Direct-Mapped Cache



The direct-mapped cache is cheaper than the fully-associative cache because it only has one tag comparator for the entire cache, as opposed to one tag comparator for every line in a fully-associative cache. However, address collisions are possible in a direct-mapped cache, which might adversely affect performance.

When a memory address is presented, the index bits will select the cache line. If the tag from that cache line is equal to the tag from the memory address, it's a cache hit. The data from that cache line will be sent to the CPU.

If there is a cache miss and we have to go to memory, the data from memory replaces the cache line indicated by the address index. If the line being replaced is valid and dirty data, the cache must write that data to memory rather than throw it out.

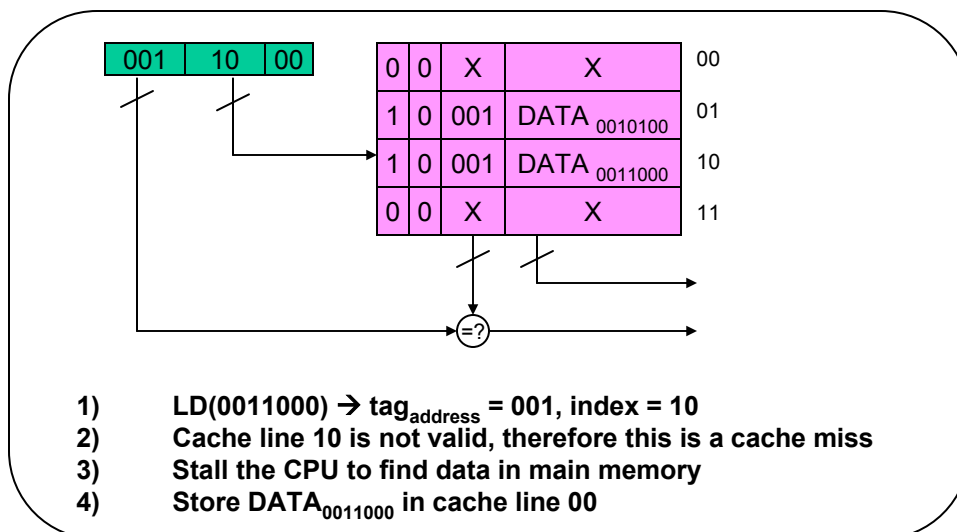
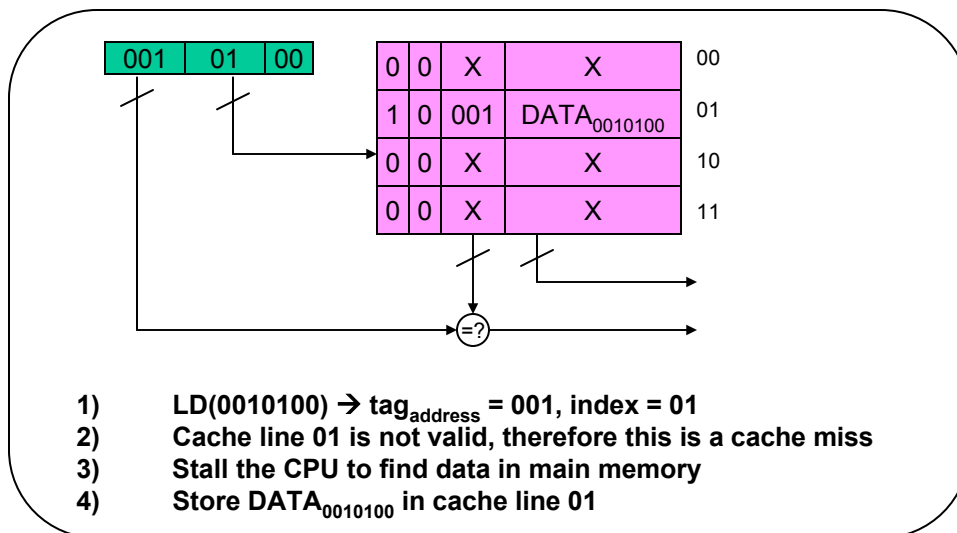
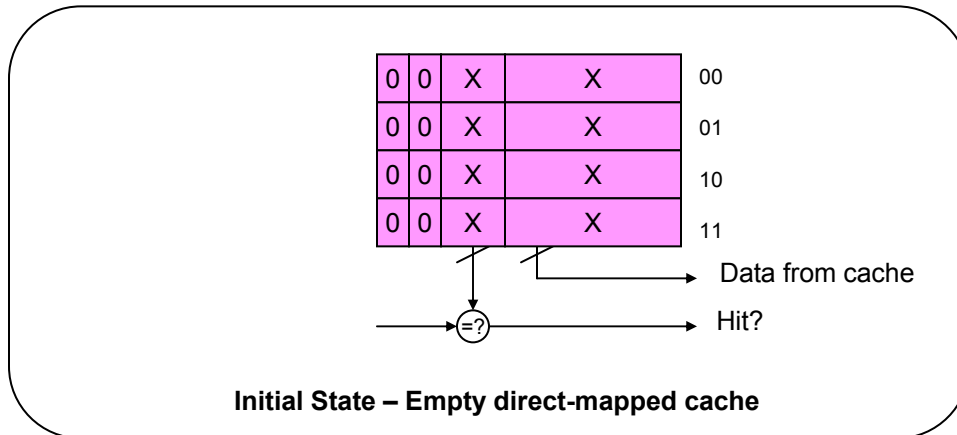
Which cache should you use? It depends on the memory access pattern. Different patterns and replacement strategies are more optimal for different caches.

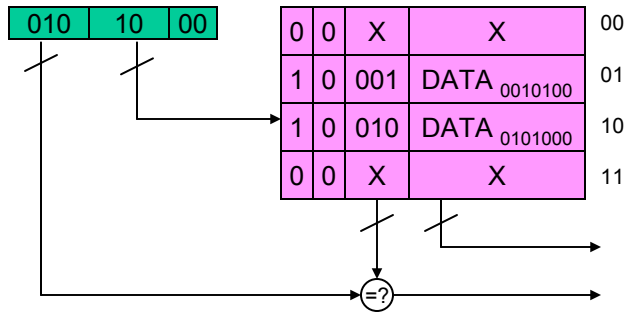
Here are some useful equations regarding direct-mapped caches:

$$\text{Capacity} = \# \text{ lines} * (1 \text{ valid bit} + S \text{ status bit} + T \text{ tag bits} + D \text{ data bits})$$

$$\# \text{ lines} = 2^{(\# \text{ index bits})}$$

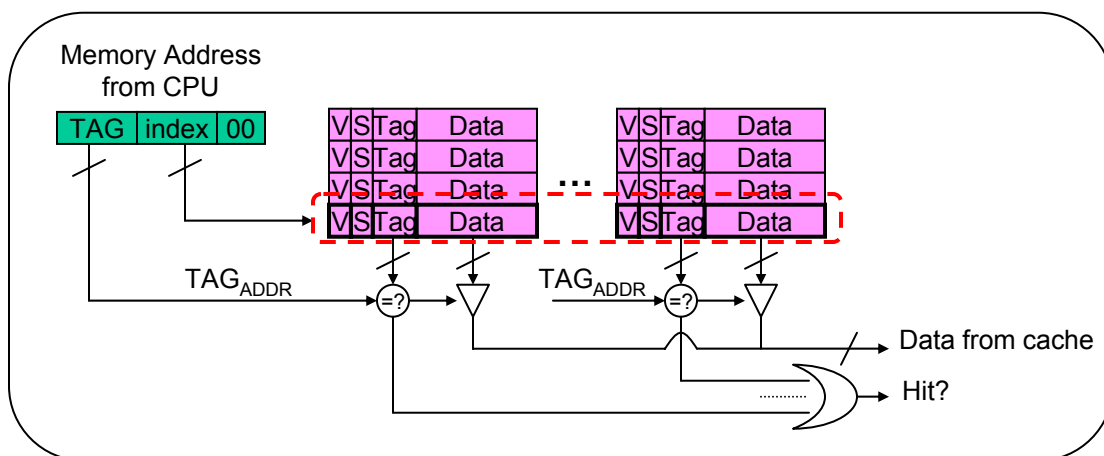
The following example illustrates four memory accesses, starting with an empty direct mapped cache. Each picture illustrates the state of the cache *after* the commentary is executed.





- 1) LD(0101000) → tag_{address} = 010, index = 10
- 2) Cache line 10 is valid
- 3) Address tag and cache line tag do not match → miss!
- 4) Stall the CPU to find data in memory
- 5) Dirty bit is zero, no need to write DATA₀₀₁₁₀₀₀ to memory
- 6) Store DATA₀₁₀₁₀₀₀ in cache line 10

N-way Set-Associative Cache



A N-way set associative cache is made up of N ways (N direct-mapped caches). Data can be stored in any one of those ways.

The set-associative cache is a tradeoff between a fully-associative and a direct-mapped cache. A set-associative cache does not have a lot of expensive tags like a fully-associative cache. A set-associative cache also cuts down on address collisions by allowing data to be stored in N cache lines rather than one cache line provided by a direct-mapped cache.

When a memory address is presented to the set-associative cache, the index bits select a cache line in all the ways. If the address tag matches any of the tags in the selected cache lines, it's a cache hit. The data from that cache line will be sent to the CPU.

If there is a cache miss and we have to go to memory, the data from memory replaces a cache line in one of the direct-mapped caches. If all the selected cache lines are valid, the cache uses a replacement policy to determine which should be replaced. Also, if the line being replaced has valid and dirty data, the cache must write that data to memory.

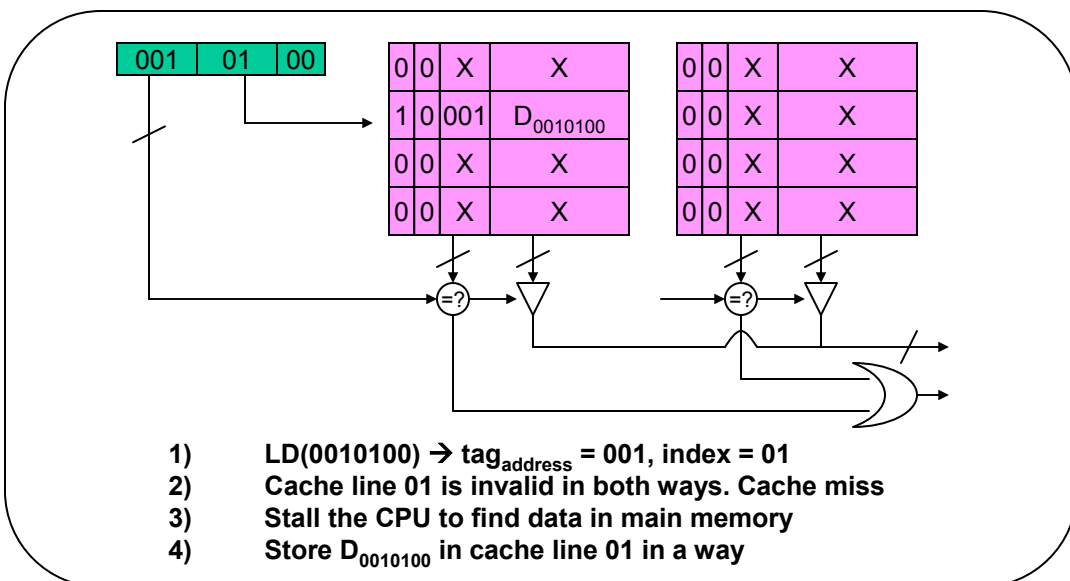
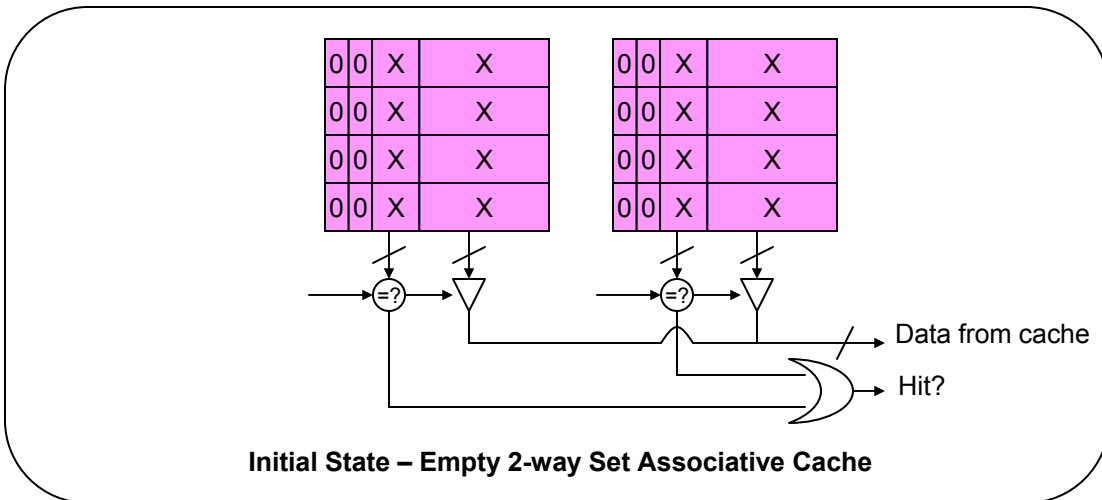
A set-associative cache has its advantages. However, that does not mean that the set-associative cache will always perform better than fully-associative or direct-mapped caches of the same capacity. As stated before, different patterns and replacement strategies are more optimal for different caches.

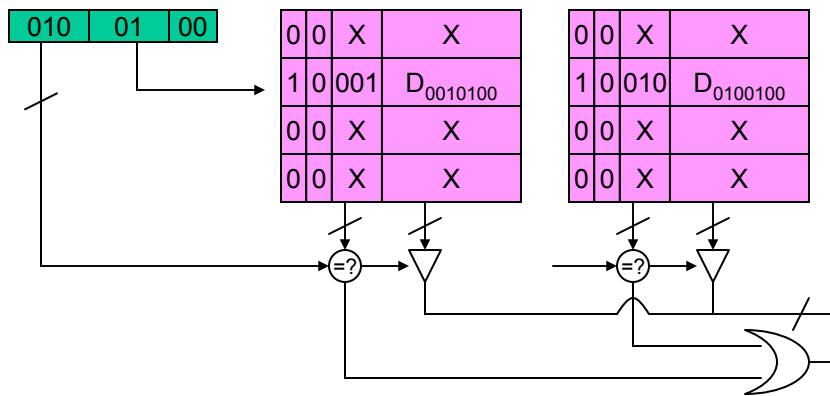
Here are some useful equations regarding set-associative caches:

$$\text{Capacity} = (\# \text{ ways}) * (\# \text{ lines/way}) * (1 \text{ valid bit} + S \text{ status bits} + T \text{ tag bits} + D \text{ data bits})$$

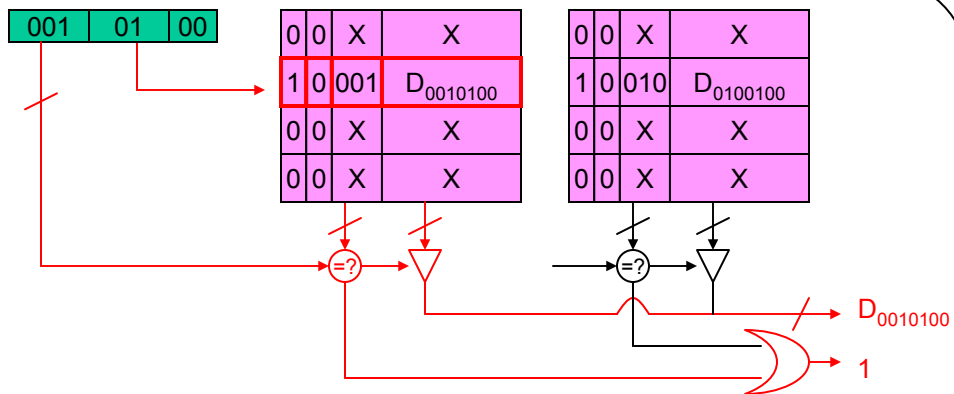
$$\# \text{ lines} = 2^{(\# \text{ index bits})}$$

The following example illustrates three memory accesses, starting with an empty 2-way set-associative cache. Each picture illustrates the state of the cache *after* the things in the commentary have executed.





- 1) LD(0100100) → tag_{address} = 010, index = 01
- 2) Cache line 01 is invalid in both ways. Cache miss
- 3) Stall the CPU to find data in main memory
- 4) Store D₀₀₁₀₁₀₀ in cache line 01 in the other way



- 1) LD(0010100) → tag_{address} = 001, index = 01
- 2) Cache Hit in first way
- 3) Return D₀₀₁₀₁₀₀ to CPU

Spatial Locality and Data Block Size

Memory accesses exhibit spatial locality – if you access memory location A then you're more likely to access locations near A in the future. Spatial locality applies to both data and instruction memory accesses.

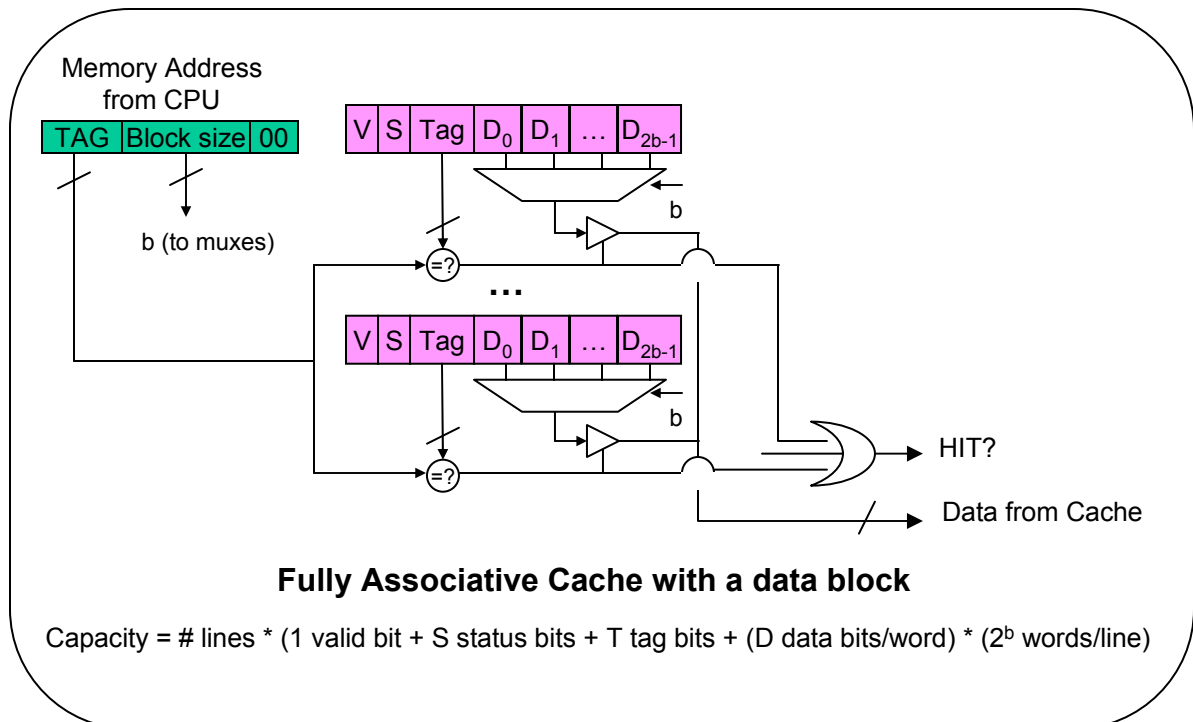
Example involving instruction accesses:

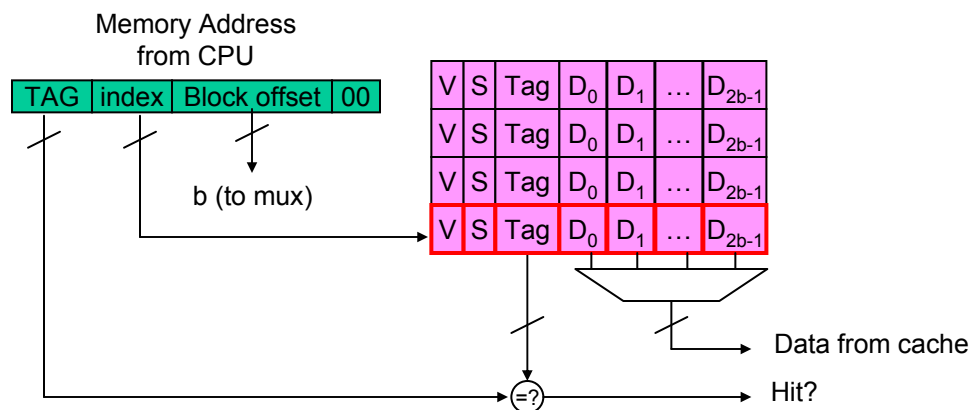
CPU instructions are usually executed sequentially through memory (with the exceptions of jumps and branches). If you fetch an instruction from memory location A, you will likely need to fetch instructions from A+4, A+8, A+12, etc. in the near future.

Caches can take advantage of spatial locality by increasing the data block size. Whenever there is a cache miss on memory location A, the cache will put a block of memory surrounding and including A into the cache.

There are some tradeoffs when it comes to block sizes (assuming a fixed cache capacity).

Advantages of a larger block size	Disadvantages of a larger block size
Miss rate decreases if the program has a large spatial locality	More address collisions → miss rate increases
	Miss penalty increases (more data is moved from memory to the cache)
	Unnecessary words fetched

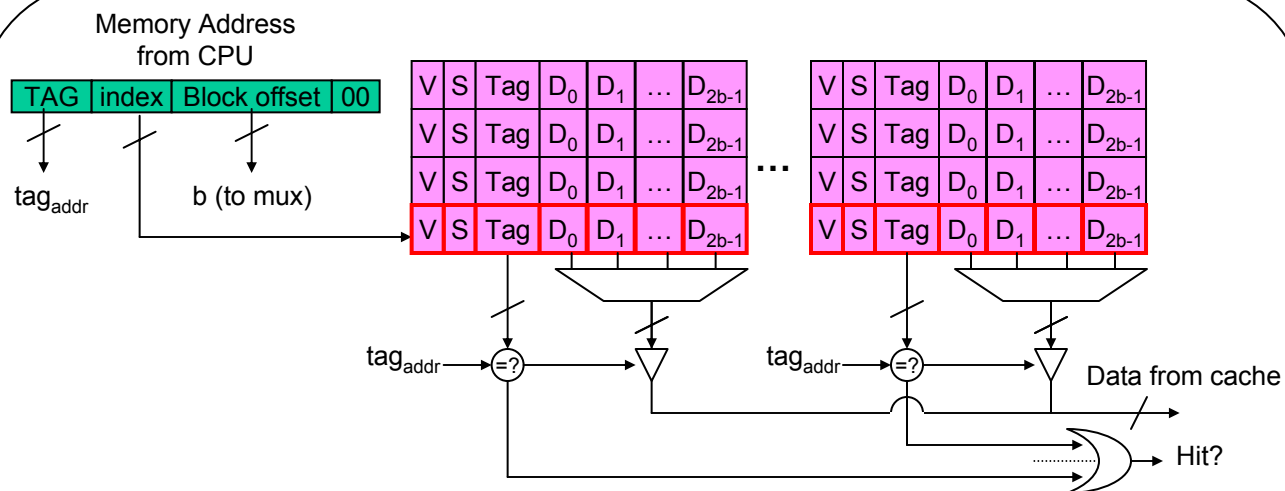




Direct Mapped Cache with a data block

Capacity = # lines * (1 valid bit + S status bits + T tag bits + (D data bits/word) (2^b words/line))

lines = $2^{(\# \text{ index bits})}$



Set Associative Cache with a data block

Capacity = (# ways) * (# lines/way) * (1 valid bit + S status bits + T tag bits + (D data bits/word) * (2^b words/line))

lines/way = $2^{(\# \text{ index bits})}$

Replacement Policies

Fully-associative and set-associative caches need a replacement policy to decide which cache line to replace when all options are valid. Three common replacement strategies include LRU (Least recently used), FIFO (first-in, first-out), and Random.

LRU replaces the least recently used location in the cache. LRU makes sense from a temporal locality point of view – if certain lines have been used more recently, they are more likely to be used in the future. LRU is rather costly, however, and is susceptible to memory patterns that cause worst-case cache performance.

FIFO replaces the oldest item in the cache. FIFO is much cheaper and easier to implement than LRU, but is also susceptible to worst-case scenario memory patterns that cause worst-case cache performance.

Random will randomly select a line to be replaced. Random generators that produce a uniform distribution are difficult to build, possibly less susceptible to worst-case scenarios.

Write Policies

When a cache line is replaced or being written, that line needs to be written to memory at some point. Three common write policies include write through, write behind, and write back.

Write through – every time the CPU writes to memory, write to both the cache and main memory. There are no dirty bits in write through. An advantage to write through is that the contents of main memory will never be stale. A disadvantage is a sacrifice in performance from stalling the CPU on every write to memory. For example, a lot of programs contain functions that repeatedly write to the same variable and then return the final result. All we need is for the final result to be stored in main memory. Storing the intermediate results in main memory might be unnecessary.

Write behind – every time the CPU writes to memory, write to the cache and the write buffer. The CPU then proceeds to the next instruction while the write buffer writes to main memory in the background. There are no dirty bits in write behind. This usually results in better performance than write-through, but it is still possible for the write buffer to unnecessarily write things to memory.

Write-back – every time the CPU writes to memory, the cache will set the dirty bit to one. On a cache miss, if the cache line being replaced has a dirty bit of one, then that line will be written to memory. Write back performs well because the cache writes to memory only when it is absolutely necessary. However, write back sacrifices some area from storing a dirty bit in each cache line.

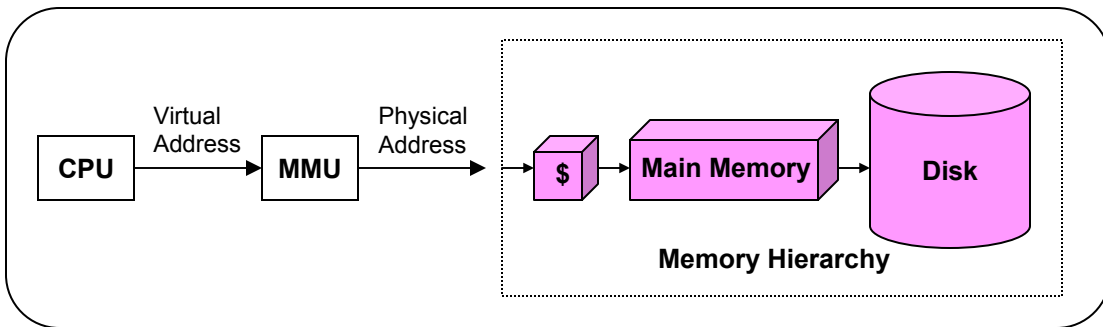
Virtual Memory

Programs that are run on the CPU commonly use *virtual* addresses rather than physical addresses when accessing memory. The CPU's virtual address is translated by the MMU (memory management unit) into a physical address that accesses the memory hierarchy.

Physical memory and virtual memory are both divided into sections called pages. Physical pages have physical page numbers (PPN) and virtual pages have virtual page numbers (VPN).

The virtual address that is inputted to the MMU contains a virtual page number and a page offset. The MMU translates the virtual page number into the appropriate physical page number via a lookup table. The MMU does not touch the page offset.

The output of the MMU is a physical address made up of a physical page number and an offset. The physical page number points to a page in memory, while the offset indicates which word to access on that page.



Using virtual memory has numerous advantages over directly addressing physical memory. Some (not all) of these advantages are discussed below.

Virtual memory separates hardware from software, thus allowing software to be written without being constrained by the size of physical memory. This allows programs to use a huge a virtual memory space that's larger than physical memory and split up the jobs between multiple computers. It also allows machines to have more physical memory than they have in virtual memory.

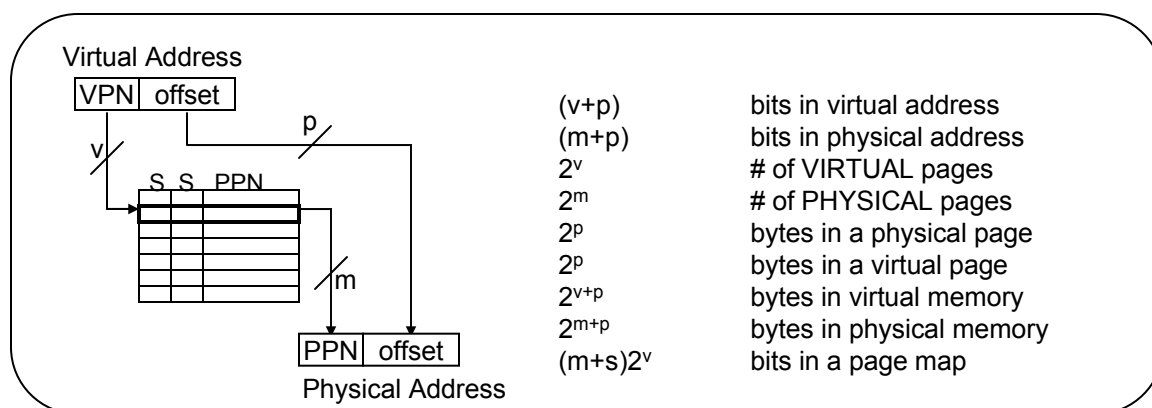
Virtual memory allows different programs to be isolated from one another, and allows the operating system to isolate itself from other software. It also allows programs and their data to be rearranged in physical memory by the operating system in a way that's transparent to the program.

Linear Page Table (one-level page map)

A page table is basically a list of PPNs somewhere in memory. The virtual address's VPN selects which PTE (page table entry) to use. The PPN in that PTE is concatenated with the offset to form a physical address.

A PTE contains optional status bits (dirty, resident, read-only, supervisor, etc.), and a PPN (physical page number).

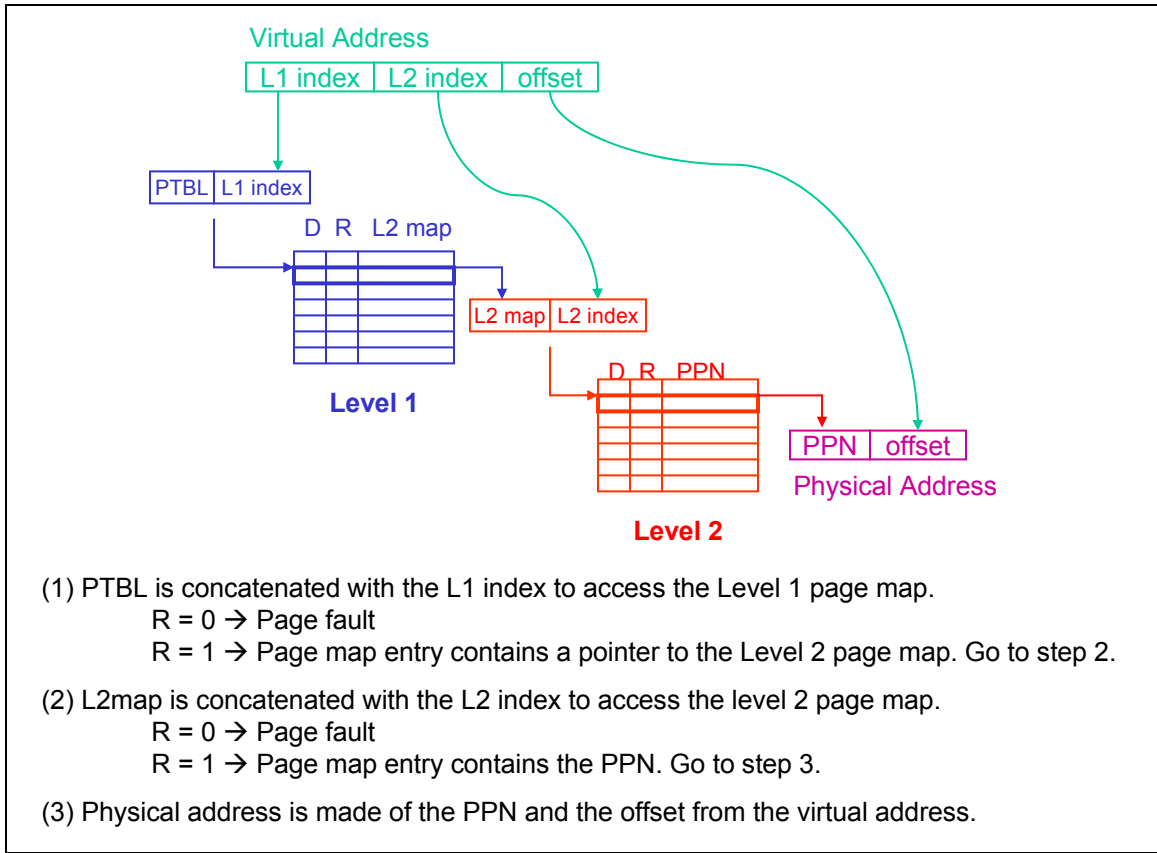
A dirty bit says the page has been changed since it was loaded from disk, and therefore needs to be written to disk when it is replaced. The resident bit is set to one for pages that are in main memory, and zero (page fault) if the page is unallocated or on disk.



Hierarchical Page Map

In a linear page table (one-level page table), we often have a situation where a lot of memory is unused, resulting in a lot of wasted space. One possible way to alleviate the waste of space is to use a hierarchical page map instead of a linear page table. The diagram below illustrates a two-level page map containing one Level 1 page table, and $2^{L1\text{index}}$ possible Level 2 page tables.

Multiple unused PTEs in a one-level page map can be equivalent to one unused Level 1 PTEs in a hierarchical page map. This significant savings in memory, however, comes at a cost in performance since page maps are usually stored in physical memory. A linear page table walk takes one memory access while a two-level walk takes two memory accesses.



Translation Look-aside Buffer (TLB)

The TLB is a small, usually fully-associative cache that contains page map entries that translate VPNs to PPNs. A TLB miss can severely hinder performance because the page table walk can result in multiple accesses to memory or a page fault.

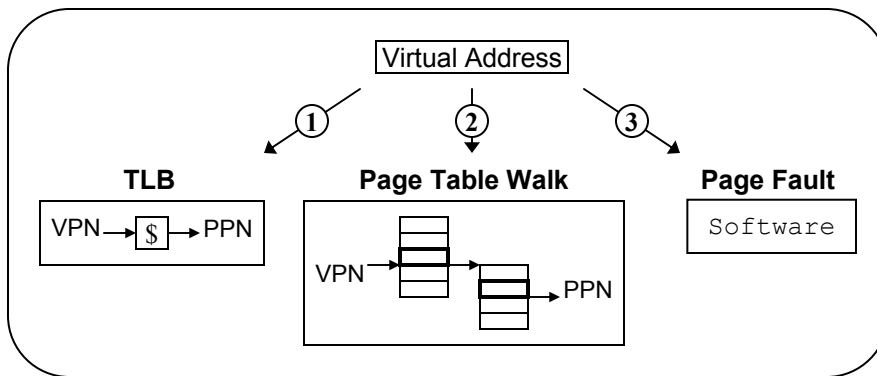
Caches – Virtual or Physical addresses

Caches that work with physical addresses might not be as advantageous if we access main memory during the page table walk. Caches that work with virtual addresses might perform better, but this brings up issues that are beyond the scope of 6.004. Take 6.823 if you want to learn about these issues.

MMU Overview

The MMU goes through the following steps when it is presented with a virtual address:

- 1) Check the TLB (Translation Look-aside Buffer). If there is a hit, send the physical address to memory. If not, go to step 2.
- 2) Page Table Walk – If the page tables and final physical page reside in memory (*not* in disk), then update the TLB and send the physical address to memory. However, if one of the page tables or the final physical page is not resident in memory, go to step 3.
- 3) Page Fault – a page table or final physical page resides in disk. Handled in software.



References

6.004 Lecture Notes from Fall 2000 and Spring 2003

6.823 Lecture Notes from Spring 2003