

ESE680-002 (ESE534): Computer Organization

Day 7: January 31, 2007
Energy and Power



Today

- Energy Tradeoffs?
- Voltage limits and leakage?
- Thermodynamics meets Information Theory
- Adiabatic Switching

- [This is an ambitious lecture]

At Issue

- Many now argue **power** will be the ultimate scaling limit
 - (not lithography, costs, ...)
- Proliferation of portable and handheld devices
 - ...battery size and life biggest issues
- Cooling, energy costs may dominate cost of electronics

What can we do about it?

$$E = \frac{1}{2} CV^2$$

$$\tau_{gd} = Q/I = (CV)/I$$

$$I_d = (\mu C_{OX}/2)(W/L)(V_{gs} - V_{TH})^2$$

Tradeoff

- $E \approx V^2$
- $\tau_{gd} \approx 1/V$
- We can trade speed for energy
- $E \times (\tau_{gd})^2 \approx \text{constant}$

Martin *et al.* *Power-Aware Computing*, Kluwer 2001
<http://caltechcstr.library.caltech.edu/308/>

Questions

- How far can this go?
 - (return to later in lecture)
- What do we do about slowdown?

Parallelism

- We have Area-Time tradeoffs
- Compensate slowdown with additional parallelism



- ...trade Area for Energy → Architectural Option

Penn ESE680-002 Spring2007 -- DeHon

7

Ideal Example

- Perhaps: 1nJ/32b Op, 10ns cycle
- Cut voltage in half
- 0.25nJ/32b Op, 20ns cycle
- Two in parallel to complete 2ops/20ns
- 75% energy reduction
 - Also 75% power reduction

Penn ESE680-002 Spring2007 -- DeHon

8

Power Density Constrained Example

- Logic Density: 1 foo-op/mm²
- Energy cost: 10nJ/foo-op @ 10GHz
- Cooling limit: 100W/cm²
- How many foo-ops/cm²/s?
 - 10nJ/mm² x 100mm²/cm²=1000nJ/cm²
 - → top speed 100MHz
 - 100M x 100 foo-ops = 10¹⁰ foo-ops/cm²/s

Penn ESE680-002 Spring2007 -- DeHon

9

Response

- How many foo-ops/cm²/s?
 - 10nJ/mm² x 100mm²/cm²=1000nJ/cm²
 - → top speed 100MHz
 - 100M x 100 foo-ops = 10¹⁰ foo-ops/cm²/s
- Power constraint won't let us run at 10GHz
 - might as well lower voltage, save energy

Penn ESE680-002 Spring2007 -- DeHon

10

What can we support?

$$E \times (t_{gd})^2 \approx \text{constant} \rightarrow 10nJ \times (100ps)^2 = E \times (t_{cycle})^2$$

$$100W / cm^2 = \left(\frac{10nJ}{\left(\frac{t_{cycle}}{100ps} \right)^2} \right) \times 100 \times \left(\frac{1}{t_{cycle}} \right)$$

Penn ESE680-002 Spring2007 -- DeHon

11

(Pushing through the Math)

$$(t_{cycle})^3 = \frac{10nJ \times 100 \times (100ps)^2}{100J / s}$$

$$t_{cycle} = \sqrt[3]{10^{-8} \times (10^{-10})^2 s^3}$$

$$t_{cycle} = 4.64 \times 10^{-10} s \approx 500ps$$

Penn ESE680-002 Spring2007 -- DeHon

12

Improved Power

- How many foo-ops/cm²/s?
 - 2GHz x 100 foo-ops = 2 × 10¹¹ foo-ops/cm²/s
 - At 5x lower voltage
 - [vs. 100M x 100 foo-ops = 10¹⁰ foo-ops/cm²/s]

How far?

Limits

- Ability to turn off the transistor
- Noise
- Parameter Variations

Sub Threshold Conduction

- To avoid leakage want I_{off} very small
- Use I_{on} for logic – determines speed
- Want I_{on}/I_{off} large

$$I_{off} = I_{VT} \times 10^{-(V_T/S)}$$

$$S = (\ln(10))\eta kT / e$$

[Frank, IBM J. R&D v46n2/3p235]

Sub Threshold Conduction

- S ≈ 90mV for single gate
- S ≈ 70mV for double gate
- 4 orders of magnitude I_{VT}/I_{off} → V_T > 280mV

$$I_{off} = I_{VT} \times 10^{-(V_T/S)}$$

$$S = (\ln(10))\eta kT / e$$

[Frank, IBM J. R&D v46n2/3p235]₁₇

ITRS2005 – High Performance

Year of Production	2005	2006	2007	2008	2009	2010	2011	2012	2013
DRAM % Pitch Area (contacts)	80	70	65	57	50	45	40	36	32
MPI/ASIC Metal 2 (M2) % Pitch Area (contacts)	90	78	68	59	52	45	40	36	32
MPI Physical Gate Length (nm)	32	28	25	22	20	18	16	14	13
<i>I_{off}</i> Physical Logic for High Performance logic (nA)	32	28	25	22	20	18	16	14	13
<i>V_{off}</i> Substrate Threshold Voltage [V]									
Exceeded Power Bulk (mV)	195	168	165	160	159	151	146	148	
UTB FD (mV)				169	168	167	170	166	167
DG (mV)							181	184	185
<i>I_{off}</i> Source Drain Subthreshold Off-State Leakage Current [fA]									
Exceeded Power Bulk (nA/μm)	0.06	0.15	0.2	0.2	0.22	0.28	0.32	0.34	
UTB FD (nA/μm)				0.17	0.16	0.22	0.32	0.39	0.39
DG (nA/μm)							0.1	0.11	0.11
<i>I_{off}</i> Effective M2O2 Drive Current [fA]									
Exceeded Power Bulk (nA/μm)	1020	1120	1200	1370	1516	1690	1896	2090	
UTB FD (nA/μm)				1486	1625	1815	2015	2207	2198
DG (nA/μm)							1889	1932	2220

ITRS2005 – Low Power

Year in Production	2005	2006	2007	2008	2009	2010	2011	2012	2013
DRAM 1s Pitch (nm) (contacted)	80	70	65	57	50	45	40	36	32
MPCU/ASIC Metal 1 (M1) 1s Pitch (nm)(contacted)	90	78	68	59	52	45	40	36	32
MPCU Physical Gate Length (nm)	22	20	20	22	20	18	16	14	13
L_{gp} Physical gate length for LOP (nm) [1]	45	37	32	28	25	22	20	18	16
V_{DD} Standby Threshold Voltage [7]									
Extended Finstr Bulk (mV)	288	303	285	274	275	226	233	231	
UTB FD (mV)							273	268	272
DG (mV)							261	255	257
$I_{D,stat}$ Source Drain Subthreshold Off-State Leakage Current [8]									
Extended Finstr Bulk (μ A/ μ m)	3.0E-03	3.0E-03	5.0E-03	5.0E-03	5.0E-03	5.0E-03	1.8E-02	2.9E-02	
UTB FD (μ A/ μ m)							8.0E-03	1.0E-02	1.0E-02
DG (μ A/ μ m)							5.0E-03	7.0E-03	7.0E-03
$I_{D,eff}$ effective NANO3 Drive Current [9]									
Extended Finstr Bulk (μ A/ μ m)	589	607	573	712	775	749	749	774	
UTB FD (μ A/ μ m)							740	765	778
DG (μ A/ μ m)							783	822	789

Penn ESE680-002 Spring2007 -- DeHon

Table 41c

19

Thermodynamics

Penn ESE680-002 Spring2007 -- DeHon

20

Lower Bound?

- Reducing entropy costs energy
- Single bit gate output
 - Set from previous value to 0 or 1
 - Reduce state space by factor of 2
 - Entropy: $\Delta S = k \times \ln(\text{before/after}) = k \times \ln 2$
 - Energy = $T \Delta S = kT \times \ln(2)$
- Naively setting a bit costs at least $kT \times \ln(2)$

Penn ESE680-002 Spring2007 -- DeHon

21

Numbers (ITRS 2005)

- $kT \times \ln(2) = 2.87 \times 10^{-21} \text{J}$ (at R.T K=300)

Year in Production	2014	2015	2016	2017	2018	2019	2020
DRAM 1s Pitch (nm) (contacted)	28	25	22	20	18	16	14
MPCU/ASIC Metal 1 (M1) 1s Pitch (nm)(contacted)	28	25	22	20	18	16	14
MPCU Physical Gate Length (nm)	11	10	9	8	7	6	6
L_{gp} Physical gate length for LOP (nm) [1]	14	13	11	10	9	8	7
V_{DD} Power Supply Voltage (V) [5]	0.6	0.6	0.5	0.5	0.5	0.5	0.5
$C_{g,stat}$ Total gate capacitance for calculation of CTI^2 [14]							
Extended Finstr Bulk (F/ μ m)	4.83E-16	5.44E-16	6.05E-16				
UTB FD (F/ μ m)	6.43E-16	6.14E-16	5.95E-16	5.24E-16	4.82E-16	4.41E-16	4.00E-16
DG (F/ μ m)							

$$W/L=3 \rightarrow W=21\text{nm}=0.021\mu\text{m}$$

$$C \approx 8 \times 10^{-18} \text{F} \approx 10^{-17} \text{F}$$

Table 41d

$$E_{op} = CV^2 = 2.5 \times 10^{-18} \text{F}$$

Penn ESE680-002 Spring2007 -- DeHon

22

Sanity Check

- $V=0.5\text{V}$
- $Q=CV=0.5 \times 10^{-17}$ coulombs
- $e=1.6 \times 10^{-19}$ coulombs
- $Q \approx 30$ electrons?
- Energy in α particle?
 - 10^5 — 10^6 electrons?

Penn ESE680-002 Spring2007 -- DeHon

23

Hmm...

- $CV^2=2.5 \times 10^{-18} \text{J}$
- 18 Billion Transistors in 2.5cm^2
 - Generous, assumes no interconnect capacitance
- $4.5 \times 10^{-8} \text{J} / 2.5\text{cm}^2 \approx 2 \times 10^{-8} \text{J/cm}^2$
- Cooling limit of @100W/cm²
- Maximum operating frequency?
- 5GHz

Penn ESE680-002 Spring2007 -- DeHon

24

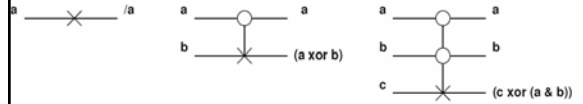
Recycling...

- Thermodynamics only says we have to dissipate energy if we discard information
- Can we compute without discarding information?
- Can we use this?

Penn ESE680-002 Spring2007 -- DeHon

25

Three Reversible Primitives



Penn ESE680-002 Spring2007 -- DeHon

26

Universal Primitives

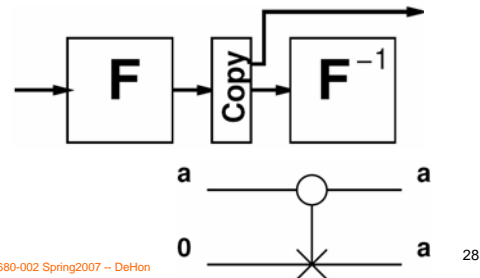
- These primitives
 - Are universal
 - Are all reversible
- If keep all the intermediates they produce
 - Discard no information
 - Can run computation in reverse

Penn ESE680-002 Spring2007 -- DeHon

27

Cleaning Up

- Can keep "erase" unwanted intermediates with reverse circuit



Penn ESE680-002 Spring2007 -- DeHon

28

Thermodynamics

- In theory, at least, thermodynamics does not demand that we dissipate any energy (power) in order to compute

Penn ESE680-002 Spring2007 -- DeHon

29

Adiabatic Switching

Penn ESE680-002 Spring2007 -- DeHon

30

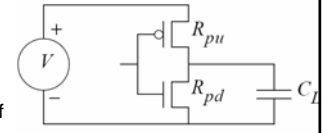
Two Observations

1. Dissipate power through on-transistor charging capacitance
2. Discard capacitor charge at end of cycle

Charge Cycle

- Charging capacitor

- $Q=CV$
- $E=QV$
- $E=CV^2$
 - Half in capacitor, half dissipated in pullup



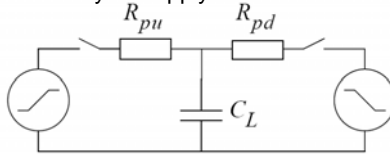
[Athas/Koller/Svenson, USC/ISI AC MOS-TR-2 1993]

Adiabatic Switching

- Current source charging:

- Ramp supplies slowly so supply constant current

- $P=I^2R$
- $E_{total}=P*T$
- $Q=IT=CV$
- $I=CV/T$
- $E_{total}=I^2R*T=(CV/T)^2R*T$
- $E_{total}=I^2R*T=(RC/T) CV^2$



Ignores leakage ...
May require large V_t

Impact of Adiabatic Switching

- $E_{total}=I^2R*T=(RC/T) CV^2$
- $RC=\tau_{gd}$
- $E_{total} \propto (\tau_{gd}/T)$
- Without reducing V
 - Can trade energy and time
- $E*T=constant$

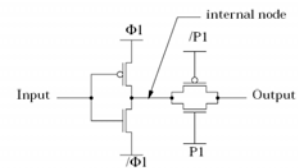
Adiabatic Discipline

- Never turn on a device with a large voltage differential across it.
- $P=\Delta V^2/R$

SCRL Inverter

- Φ 's, nodes, at $V_{dd}/2$
- P1 at ground

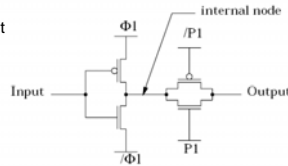
- Slowly turn on P1
- Slow split Φ 's
- Slow turn off P1's
- Slow return Φ 's to $V_{dd}/2$



[Younis/Knight ISLPED(?) 1994]

SCRL Inverter

- Basic operation
 - Set inputs
 - Split rails to compute output adiabatically
 - Isolate output
 - Bring rails back together
- Have transferred logic to output
- Still need to worry about resetting output adiabatically

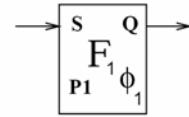


37

Penn ESE680-002 Spring2007 -- DeHon

SCRL NAND

- Same basic idea works for nand gate
 - Set inputs
 - Adiabatically switch output
 - Isolate output
 - Reset power rails

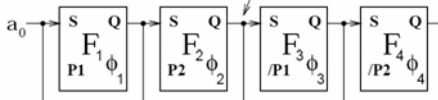


38

Penn ESE680-002 Spring2007 -- DeHon

SCRL Cascade

- Cascade like domino logic
 - Compute phase 1
 - Compute phase 2 from phase 1...
- How do we restore the output?

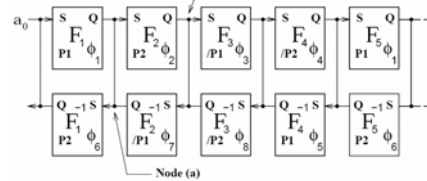


39

Penn ESE680-002 Spring2007 -- DeHon

SCRL Pipeline

- We must **uncompute** the logic
 - Forward gates compute output
 - Reverse gate restore to $V_{dd}/2$



40

Penn ESE680-002 Spring2007 -- DeHon

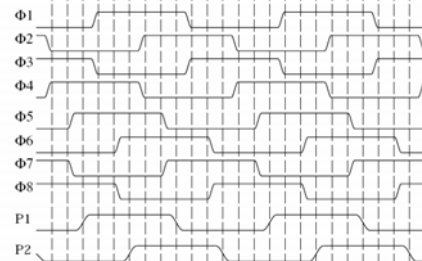
SCRL Pipeline

- P1 high (F1 on; F1 inverse off)
 - Φ_1 split: $a = F_1(a_0)$
 - Φ_2 split: $b = F_2(F_1(a_0))$
- $F_2^{-1}(F_2(F_1(a_0))) = a$
- P1 low – now F_2^{-1} drives a
- F1 restore by Φ_1 converge
- ...restore F2
- Use F_2^{-1} to restore a to $V_{dd}/2$ adiabatically

41

Penn ESE680-002 Spring2007 -- DeHon

SCRL Rail Timing



42

Penn ESE680-002 Spring2007 -- DeHon

SCRL

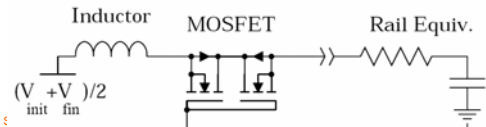
- Requires Reversible Gates to uncompute each intermediate
- All switching (except IO) is adiabatic
- Can, in principle, compute at any energy

Penn ESE680-002 Spring2007 -- DeHon

43

Trickiness

- Generating the ramped clock rails
- Use LC circuits
- Need high-Q resonators
- Making this efficient is key to **practical** implementation
 - Some claim not possible in practice



Penn ESE680-002

Big Ideas

- Can trade time for energy
 - ...area for energy
- Noise and subthreshold conduction limit voltage scaling
- Thermodynamically admissible to compute without dissipating energy
- Adiabatic switching alternative to voltage scaling
- Can base CMOS logic on these observations

Penn ESE680-002 Spring2007 -- DeHon

45