



**FAMU-FSU College of Engineering**



# **Reversible Computing**

## **A Requirement for Extreme Supercomputing**

**Dr. Michael P. Frank, Assistant Professor  
Dept. of Electrical & Computer Eng.  
FAMU-FSU College of Engineering**

**ECE Department Graduate Seminar  
Thursday, September 2, 2004**

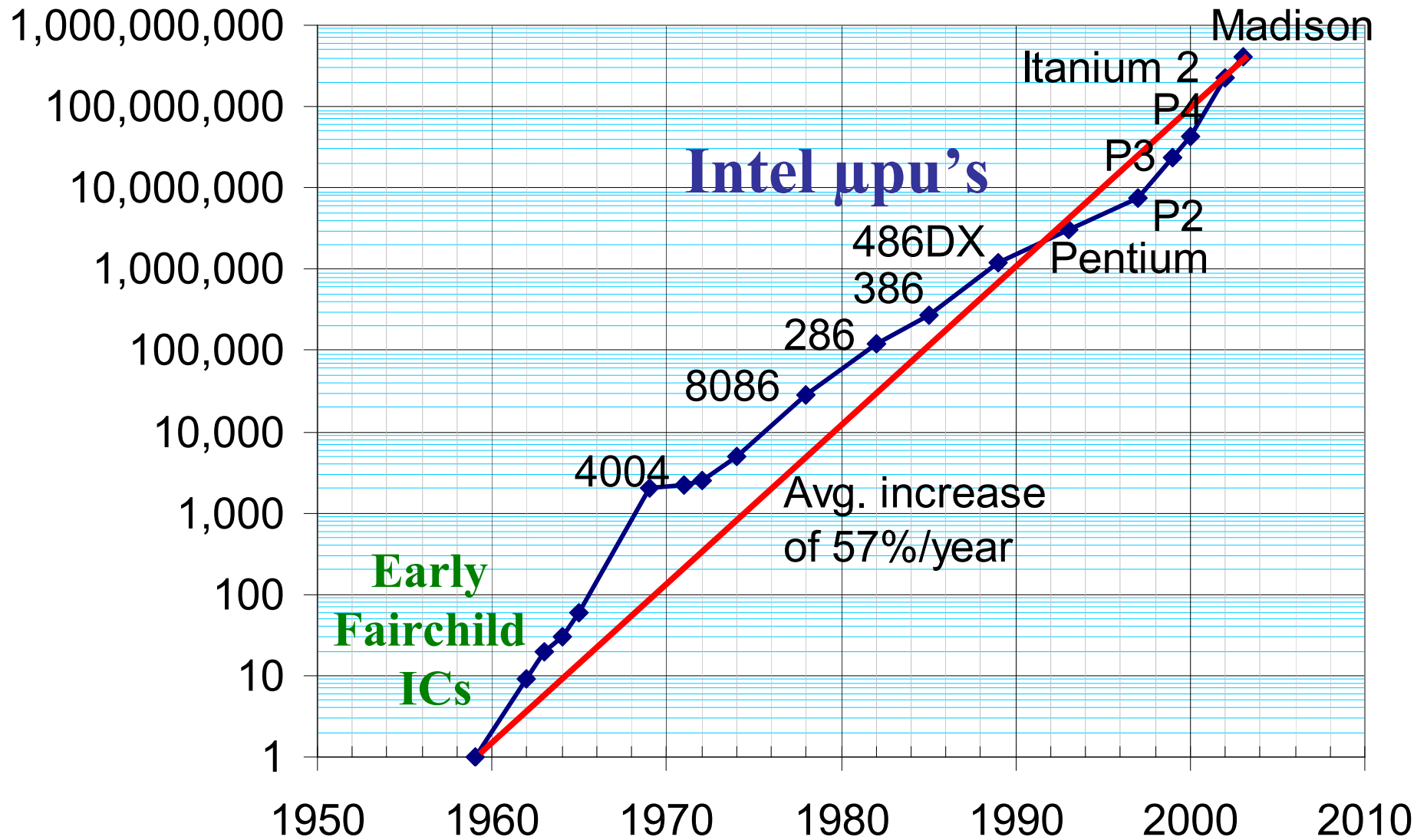


# Abstract

- The energy dissipated per switching event directly limits any digital system's performance per unit of power consumption.
  - *E.g.*, typical logic node switching energy today =  $\sim 0.1$  fJ.
    - A 1 MW machine could do "only" 100 PFLOPs. (assuming  $10^5$  logic ops/FLOP)
    - A 1 ZFLOPs machine in today's tech. would require at least 10 GW!
      - This is the approximate electrical power consumption of Norway!
- Traditional "irreversible" switching mechanisms are subject to a relatively high minimum energy dissipation per signal transition event.
  - The practical limit for irreversible CMOS may be only  $\sim 1$  order of magnitude better than today's technology.
    - And further, *any possible* irreversible technology is at best only  $\sim 2-4$  orders of magnitude better than today's!
      - *E.g.*, 1 ZFLOPs, terrestrially  $\rightarrow$  at least  $\sim 40$  MW (non-adiabatic)
- Circumventing all these bounds will require moving to increasingly *reversible* switching mechanisms and logic styles...
  - With long-term implications for computer architecture, programming languages, and algorithm design...
- In this talk, we survey reversible computing principles.
  - We argue: Reversible computing needs to be more aggressively explored!



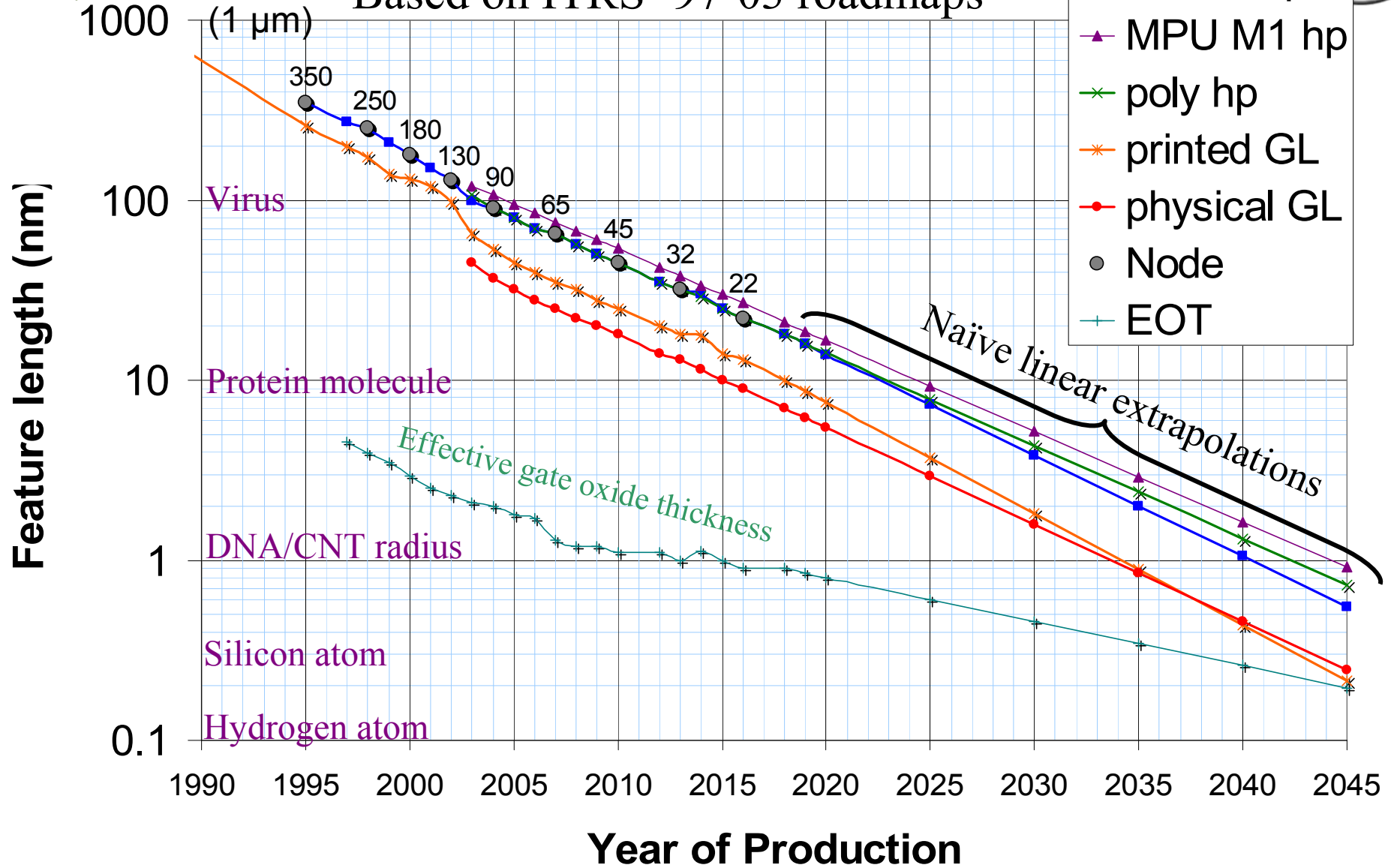
# Moore's Law (Devices/IC)





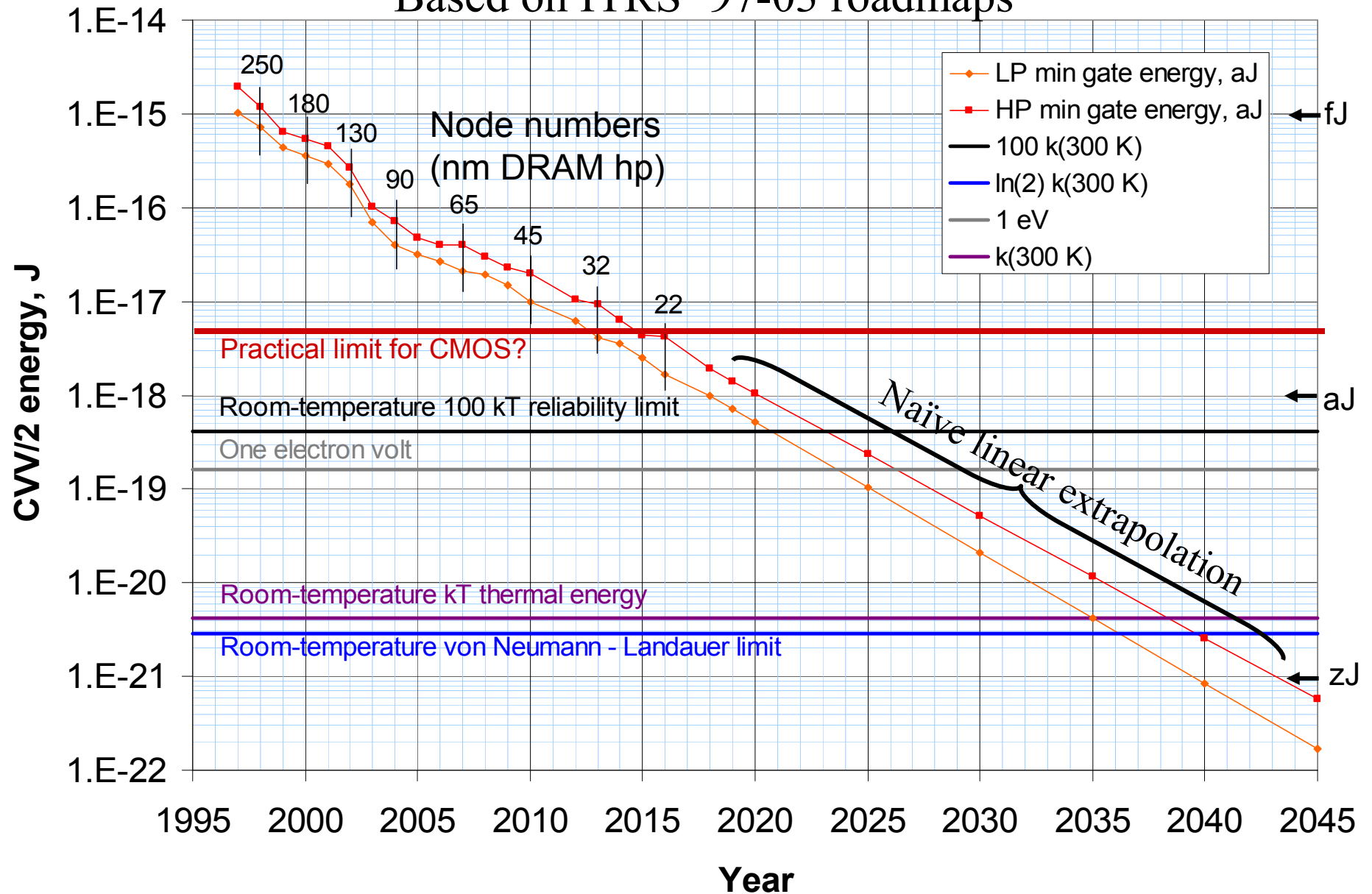
# Device Size Scaling Trends

Based on ITRS '97-03 roadmaps



# Trend of Min. Transistor Switching Energy

Based on ITRS '97-03 roadmaps





# Important Energy Limits

- Near-term leakage-based limit for MOSFETs:
  - May be  $\sim 5$  aJ, roughly  $10\times$  lower than today.
    - $10\times$  faster machines,  $\sim 4$ - $8$  years left on the clock
- Reliability-based limit on bit energies:
  - Roughly  $100 kT \approx 400$  zJ,  $\sim 100\times$  below now.
    - $100\times$  faster machines,  $\sim 8$ - $15$  years to go...
- Landauer limit on dissipation per bit erasure:
  - About  $0.7 kT \approx 3$  zJ,  $\sim 10,000\times$  below today.
    - $10,000\times$  faster machines,  $\sim 15$ - $30$  years left...
- No limit is known for reversible computing...
  - We need to investigate this alternative further.



# FET Energy Limit

- A practical limit for all transistors based on the *field effect* principle.
  - It's probably not an absolutely unavoidable, fundamental limit.
    - However, it is probably the biggest barrier to further transistor scaling today.
- The limit arises from the following chain of considerations:
  - We require reduced energy dissipation per logic operation.
  - Want small  $\frac{1}{2}CV^2$  logic node energy (normally dissipated when switching)
  - Want small node capacitance  $C$  → small transistor size (also for speed)
  - Need to lower switching voltage  $V$ , due to many factors:
    - Gate oxide breakdown, punch-through, also helps reduce  $CV^2$ .
  - Reduced on-off ratio  $R_{on/off} = I_{on}/I_{off} < e^{Vq/kT}$  (at room temperature)
    - Comes from Boltzmann (or Fermi-Dirac) distrib. of state occupancies near equil.
      - Independent of materials! (Carbon nanotubes, nanowires, molecules, etc.)
  - Increased off-state current  $I_{off}$  and power  $I_{off}V$ , given high-performance  $I_{on}$ .
  - Also, increased per-area leakage current due to gate oxide tunneling, etc.
  - Previous two both *increase* total per-device power consumption floor
    - Adds to total energy dissipated per logic gate, per clock cycle
- Eventually, the extra power dissipation from leakage overwhelms the power/performance reductions that we would gain by reducing  $CV^2$ !
  - Beyond this point, further transistor scaling hurts us, rather than helping.
    - Transistor scaling then halts, for all practical purposes!



# Mitigating MOSFET Limits

- Reduce the portion of the  $\frac{1}{2}CV^2$  node energy that gets *dissipated*
  - Reversible computing with adiabatic circuits does this
- Reduce parasitic capacitances that contribute to logic node's  $C$ 
  - via silicon-on-insulator (SOI) devices, low-k field dielectric materials, *etc.*
- Use high-k gate dielectric materials →
  - Allows gate dielectrics to be thicker for a given capacitance/area
  - Reduces tunneling leakage current through gate dielectric. Also:
  - Avoids gate oxide breakdown → allows higher  $V$ 
    - indirectly helps reduce off-state conduction.
- Use multi-gate structures (FinFET, surround-gate, *etc.*) to
  - reduce subthreshold slope  $s = V/(\log R_{on/off})$  to approach theoretical optimum,
    - $s = T/q = (kT/q \ln 10)/\text{decade} = 60 \text{ mV/decade}$
- Use multi-threshold devices & power-management architectures to turn off inactive devices to suppress leakage in unused portions of the chip
  - The remaining leakage in the active logic is still a big problem, however...
- Lower operating temperature to increase  $Vq/kT$  and thus  $I_{DS}$  on-off ratio?
  - May also lead to problems with carrier concentration, cooling costs, *etc.*
  - Conflicts with the high generalized temperature of high-frequency logic signals
- Consider devices using non-field-effect based switching principles:
  - Y-branch, quantum-dot, spintronic, superconducting, (electro)mechanical, *etc.*





# Reliability-Based Limit

- A limit on signal (bit) energy.
- Applies to any mechanism for *storing* a bit whose operation is based on the *latching* principle, namely:
  - We have some physical entity whose state (e.g. its location) encodes a bit.
    - E.g., could be a packet of electrons, or a mechanical rod
  - If the bit is 1, the entity gets “pushed into” a particular state and held there by a potential energy difference (between there and not-there) of  $E$ .
    - The entity sits in there at thermal equilibrium with its environment.
  - A potential energy barrier is then raised in between the states, to “latch” the entity into place (if present).
    - A transistor is turned off, or a mechanical latching mechanism is locked down
- The Boltzmann distribution implies that  $E > T \log N = kT \ln N$ , in order for the probability of incorrect storage to be less than  $1/N$ .
  - For electrons (fermions), we must use the Fermi-Dirac distribution instead...
    - But this gives virtually identical results for large  $N$ .
- When erasing a stored bit, typically we would dissipate the energy  $E$ .
  - However, this limit might be avoidable via special level-matching, quasi-adiabatic erasure mechanisms, or non-equilibrium bit storage mechanisms.



# Numerical Example

- **Example:** Reliability factor of  $N=10^{27}$  (e.g., 1 error in a  $10^9$  gate processor running for ~3 years at 10 GHz)
  - The entropy associated with the per-op error probability is then:  
 $\log 10^{27} = 27 \log 10 = 27 k_B \ln 10 \approx 62 k_B = 8.6 \times 10^{-22} \text{ J/K}$
  - Heat that must be output to a room- $T$  (300 K) environment:  
 $k_B (300 \text{ K}) \ln 10^{27} = 2.6 \times 10^{-19} \text{ J}$  (or 260 zJ, or 1.6 eV)
    - Sounds small, but...
  - If each gate dumped this energy @ a frequency of 10 GHz,
    - the total power dissipated by an entire  $10^9$ -gate processor is 26 W.
    - Could have at most 4 such processors within a 100 W power budget!
  - Maximum performance:  $4 \times 10^{20}$  gate-cycles/sec.
    - or 4 PFLOPS, if processors require ~100,000 logic ops on average to carry out 1 standard (double-precision) floating-point op
      - a fairly typical figure for today's well-optimized floating-point units
    - Typical COTS microprocessors today have ~100× additional overhead,
      - Leading to 40 TFLOPS max performance if using these same architectures
        - » A 40-TFLOP supercomputer (e.g. Blue Gene/L) burns ~200 kW today
        - » Only 2,000× above the reliability-based limit!



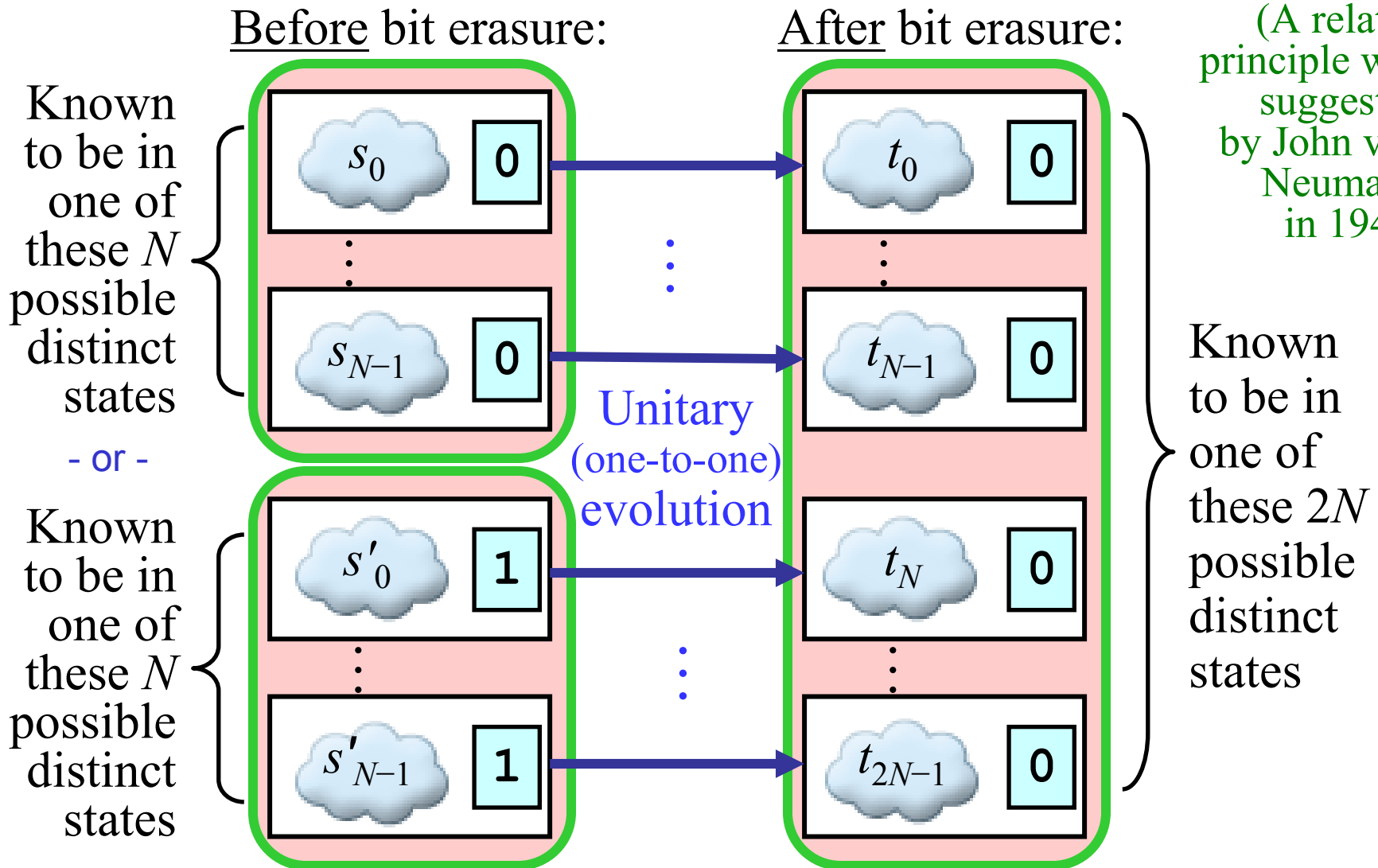
# Von Neumann / Landauer (VNL) bound for bit erasure



- The von Neumann-Landauer (VNL) lower bound for energy dissipation from bit erasure:
  - “Oblivious” erasure/overwriting of a known logical bit moves the information that it previously contained to the environment → The information becomes entropy.
    - Leads to fundamental limit of  $kT \ln 2$  for oblivious erasure.
  - This particular limit could *only* possibly be avoidable through reversible computing.
    - Reversible computing “de-computes” unwanted bits, rather than obviously erasing them!
      - This enables the signal energy to be preserved for later re-use, rather than dissipated.



# Rolf Landauer's principle (IBM Research, 1961): The minimum energy cost of oblivious bit erasure



Increase in entropy:  $\Delta S = \log 2 = k \ln 2$ . Energy dissipated to heat:  $T\Delta S = kT \ln 2$



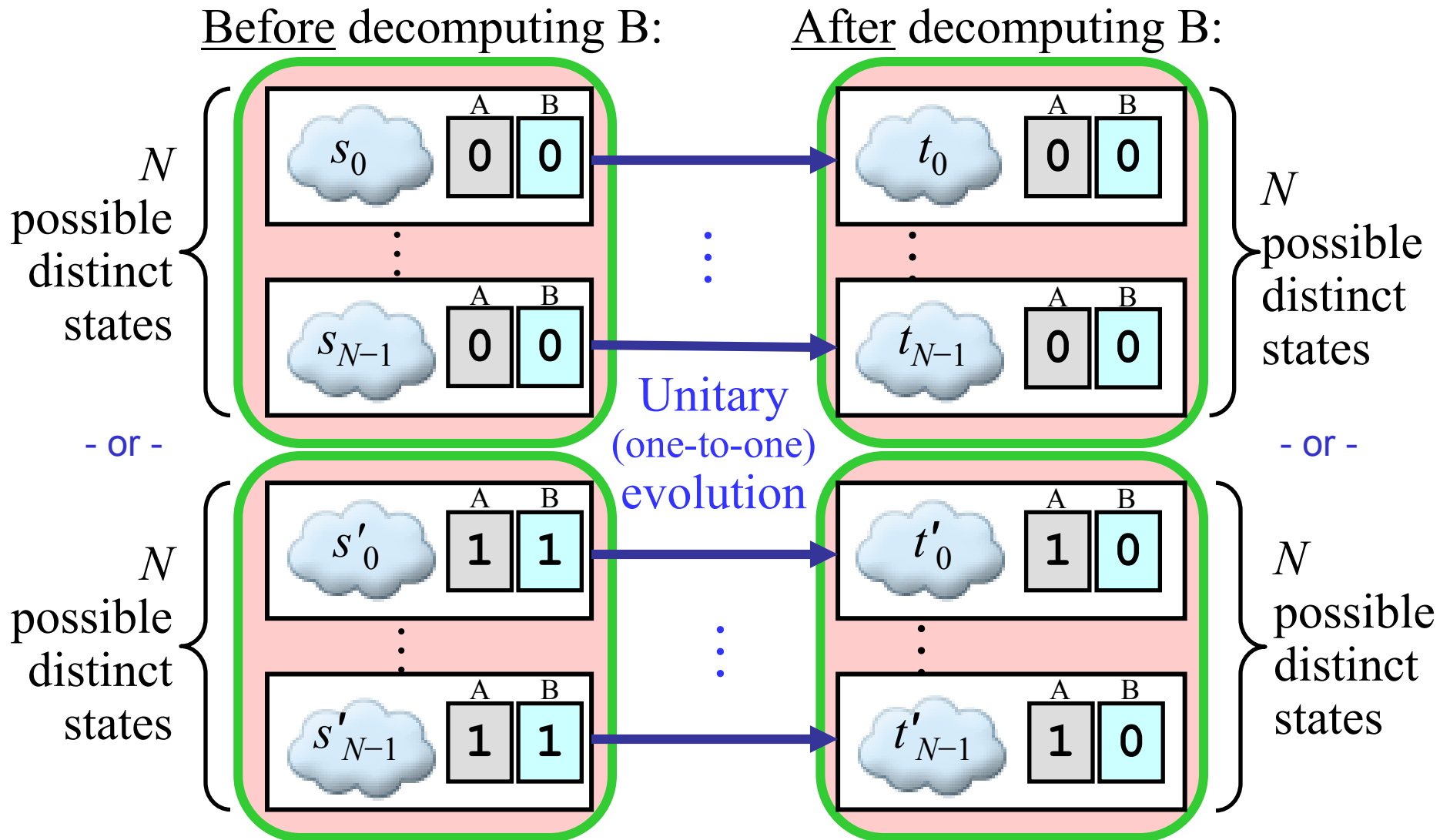
# Reversible Computing



- A *reversible* digital logic operation is:
  - Any operation that performs an invertible (one-to-one) transformation of the device's local digital state space.
    - Or at least, of that subset of states that are actually used in a design.
- Landauer's principle only limits the energy dissipation of ordinary *irreversible* (many-to-one) logic operations.
  - Reversible logic operations could dissipate much less energy,
    - Since they can be implemented in a thermodynamically reversible way.
- In 1973, Charles Bennett (IBM Research) showed how any desired computation can in fact be performed using *only* reversible logic operations (with essentially no bit erasure).
  - This opened up the possibility of a vastly more energy-efficient alternative paradigm for digital computation.
- After 30 years of (sporadic) research, this idea is finally approaching the realm of practical implementability...
  - Making it happen is the goal of the RevComp project.



# Non-oblivious “erasure” (by *decomputing* known bits) avoids the von Neumann–Landauer bound



Increase in entropy:  $\Delta S \rightarrow 0$ . Energy dissipated to heat:  $T\Delta S \rightarrow 0$



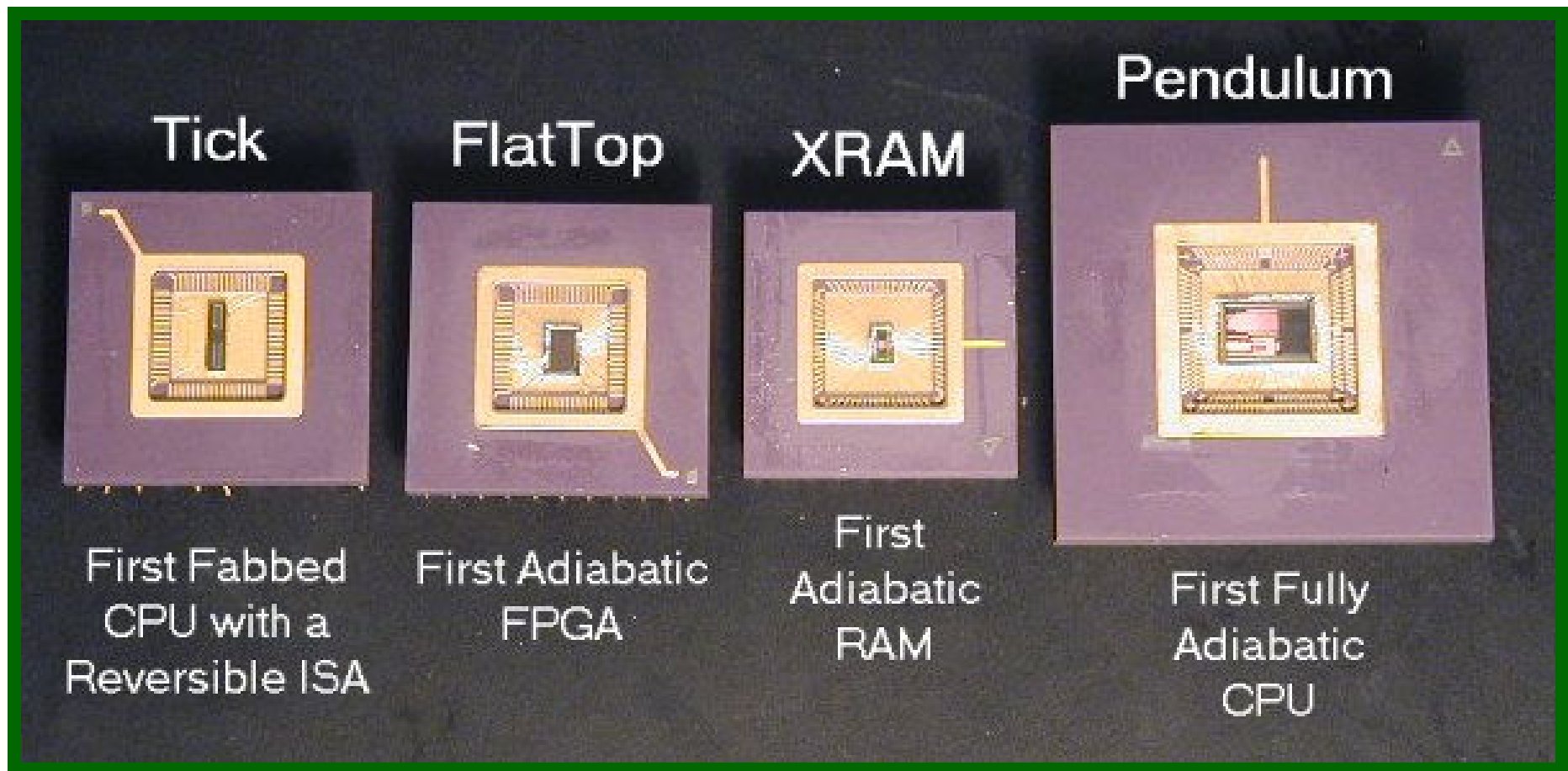
# Adiabatic Circuits

- Reversible logic can be implemented today using fairly ordinary voltage-coded CMOS VLSI circuits.
  - With a few changes to the logic-gate/circuit architecture.
- We avoid dissipating most of the circuit node energy when switching, by transferring charges in a nearly *adiabatic* (literally, “without flow of heat”) fashion.
  - *I.e.*, asymptotically thermodynamically reversible.
    - In the limit, as various low-level technology parameters are scaled.
- There are many designs for purported “adiabatic” circuits in the literature, but most of them contain fatal design flaws and are not truly adiabatic.
  - Many past designers are unaware of (or accidentally failed to meet) all the requirements for true thermodynamic reversibility.



# Reversible &/or Adiabatic VLSI Chips Designed @ MIT, 1996-1999

By Frank and other then-students in the MIT Reversible Computing group,  
under CS/AI lab members Tom Knight and Norm Margolus.







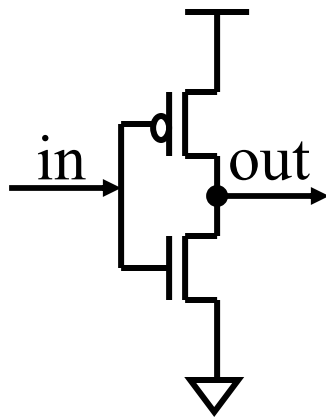
# Conventional Logic is Irreversible

Even a simple NOT gate, as it's traditionally implemented!

- Here's what all of today's logic gates (including NOT) do continually, *i.e.*, every time their input changes:
  - They overwrite previous output with a function of their input.
  - Performs many-to-one transformation of local digital state!
  - $\therefore$  required to dissipate  $\gtrsim kT$  on average, by Landauer principle
  - Incurs  $\frac{1}{2}CV^2$  energy dissipation when the output changes.

## Example:

### Static CMOS Inverter:



### Inverter transition table:

<i>Just before transition:</i>		<i>After transition:</i>	
<u>in</u>	<u>out</u>	<u>in</u>	<u>out</u>
0	0	0	1
0	1	0	1
1	0	1	0
1	1	1	0

Red arrows point from the first and fourth rows of the 'Just before' column to the first and fourth rows of the 'After' column, respectively. Green arrows point from the second and third rows of the 'Just before' column to the second and third rows of the 'After' column, respectively.

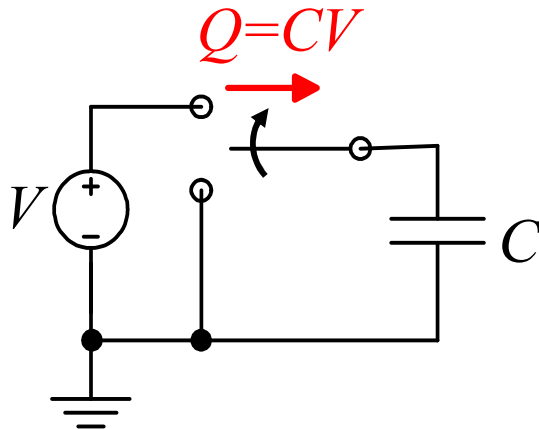


# Conventional vs. Adiabatic Charging

For charging a capacitive load  $C$  through a voltage swing  $V$

- **Conventional charging:**

- Constant voltage source:

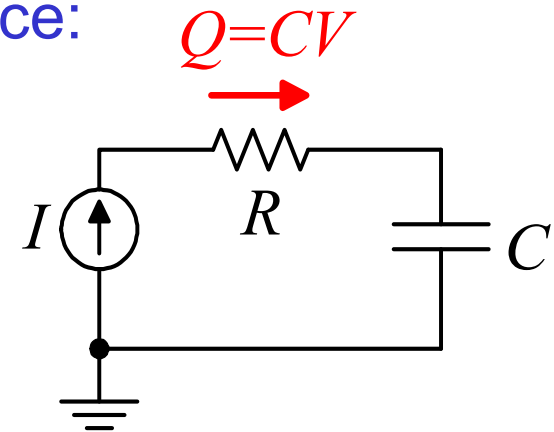


- Energy dissipated:

$$E_{\text{diss}} = \frac{1}{2} CV^2$$

- **Ideal adiabatic charging:**

- Constant current source:



- Energy dissipated:

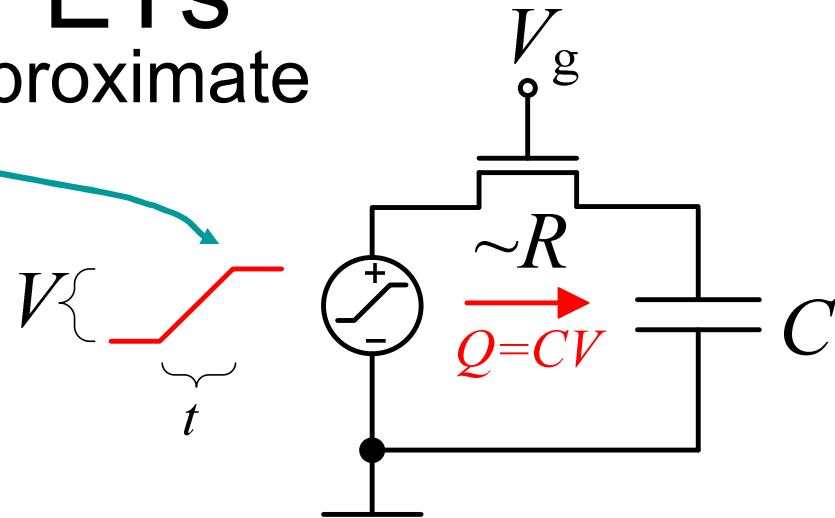
$$E_{\text{diss}} = I^2 R t = \frac{Q^2 R}{t} = CV^2 \frac{RC}{t}$$

**Note:** Adiabatic beats conventional by advantage factor  $A = t/2RC$ .



# Adiabatic Switching with MOSFETs

- Use a voltage ramp to approximate an ideal current source.
- Switch *conditionally*, if MOSFET gate voltage  $V_g > V + V_T$  during ramp.
- Can discharge the load later using a similar ramp.
  - Either through the same path, or a different path.



$$t \gg RC \Rightarrow E_{\text{diss}} \rightarrow CV^2 \frac{RC}{t}$$

$$t \ll RC \Rightarrow E_{\text{diss}} \rightarrow \frac{1}{2} CV^2$$

Exact formula:  
 $E_{\text{diss}} = s[1 + s(e^{-1/s} - 1)]CV^2$   
 given *speed fraction*  
 $s \equiv RC/t$



# Requirements for True Adiabatic Logic in Voltage-coded, FET-based circuits

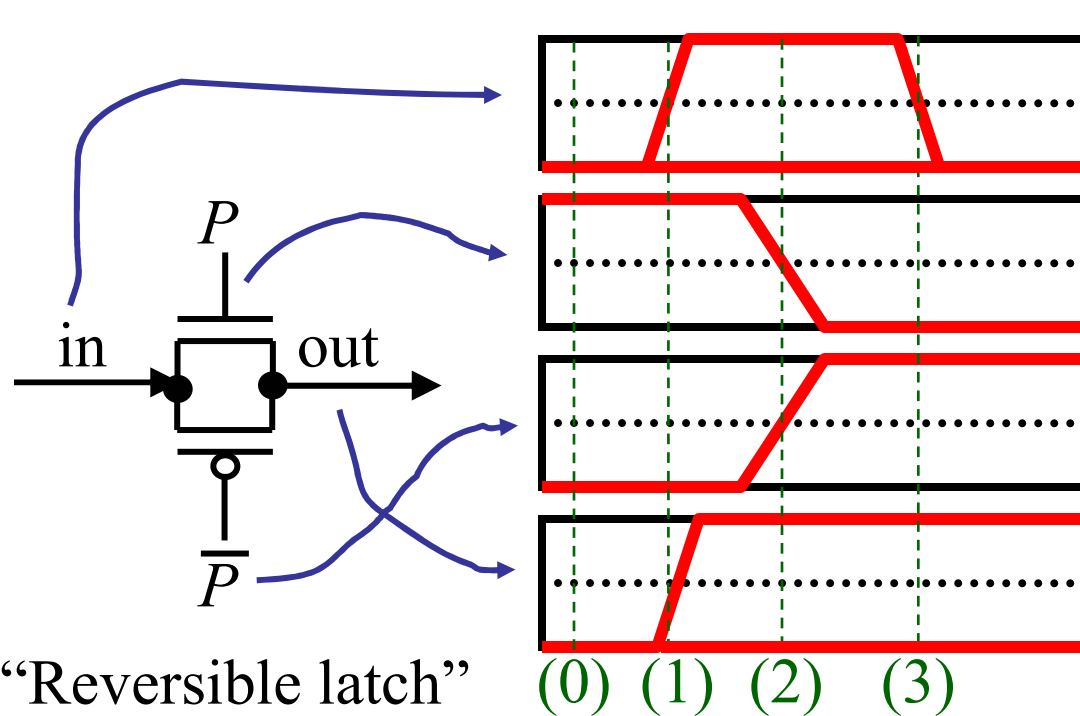
- Avoid passing current through diodes.
  - Crossing the “diode drop” leads to irreducible dissipation.
- Follow a “dry switching” discipline (in the relay lingo):
  - Never turn on a transistor when  $V_{DS} \neq 0$ .
  - Never turn off a transistor when  $I_{DS} \neq 0$ .
- Together these rules imply:
  - The logic design must be logically reversible
    - There is no way to erase information under these rules!
  - Transitions must be driven by a quasi-trapezoidal waveform
    - It must be generated resonantly, with high  $Q$
- Of course, leakage power must also be kept manageable.
  - Because of this, the optimal design point will not necessarily use the smallest devices that can ever be manufactured!
    - Since the smallest devices may have insoluble problems with leakage.

Important  
but often  
neglected!



# A Simple Reversible CMOS Latch

- Uses a single standard CMOS *transmission gate* (T-gate).
- Sequence of operation:
  - (0) input level initially tied to latch 'contents' (output);
  - (1) input changes gradually → output follows closely;
  - (2) latch closes, charge is stored dynamically (node floats);
  - (3) afterwards, the input signal can be removed.



Before input:		Input arrived:		Input removed:	
<u>in</u>	<u>out</u>	<u>in</u>	<u>out</u>	<u>in</u>	<u>out</u>
0	0	0	0	0	0
		1	1	0	1

- Later, we can reversibly “unlatch” the data with an exactly time-reversed sequence of steps.

“Reversible latch”



# 2LAL: 2-level Adiabatic Logic

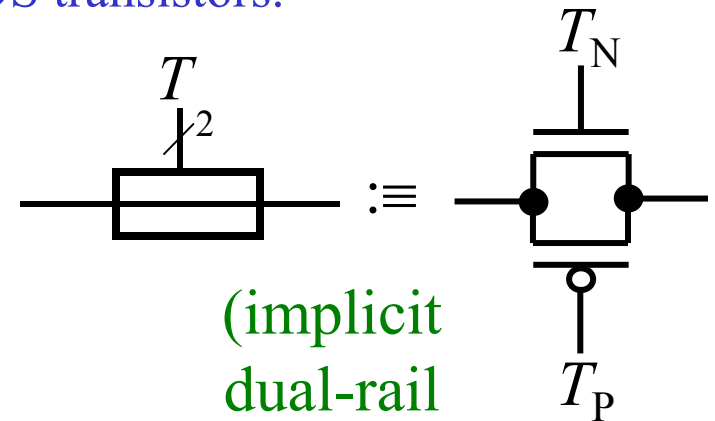
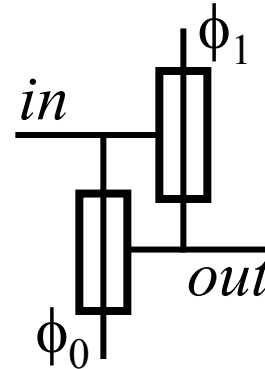


A pipelined fully-adiabatic logic invented at UF (Spring 2000), implementable using ordinary CMOS transistors.

- Use simplified T-gate symbol:
- Basic buffer element:

– cross-coupled T-gates:

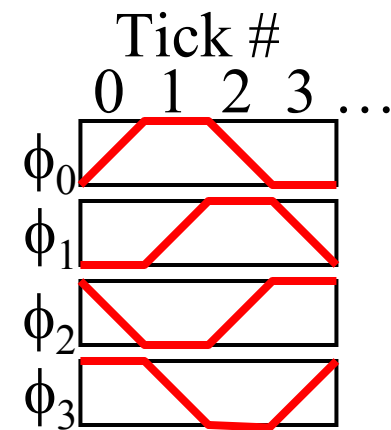
- need 8 transistors to buffer 1 dual-rail signal



(implicit dual-rail encoding everywhere)

- Only 4 timing signals  $\phi_{0-3}$  are needed. Only 4 ticks per cycle:

- $\phi_i$  rises during ticks  $t \equiv i \pmod{4}$
- $\phi_i$  falls during ticks  $t \equiv i+2 \pmod{4}$



Animation:

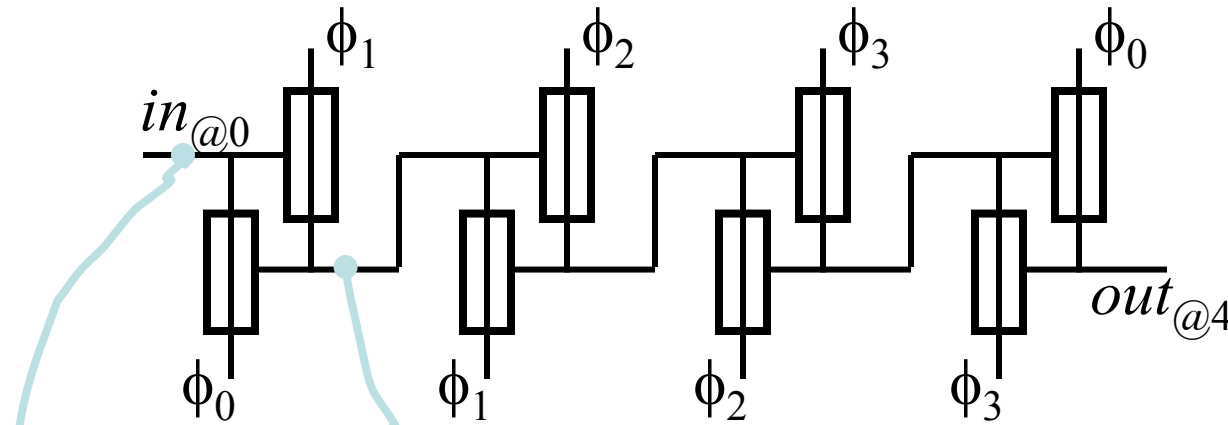




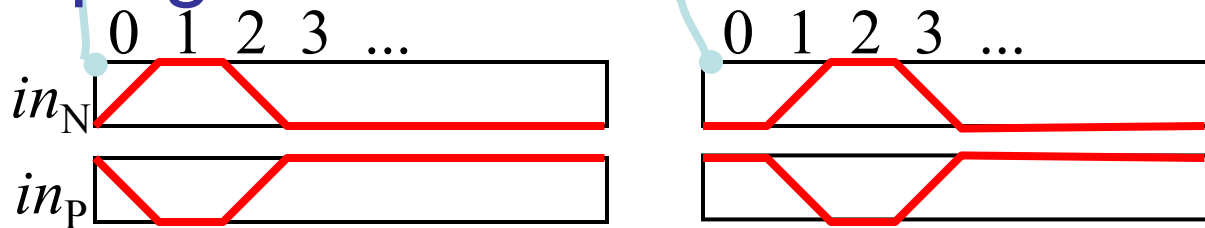
# 2LAL Shift Register Structure

- 1-tick delay per logic stage:

Animation:



- Logic pulse timing and signal propagation:

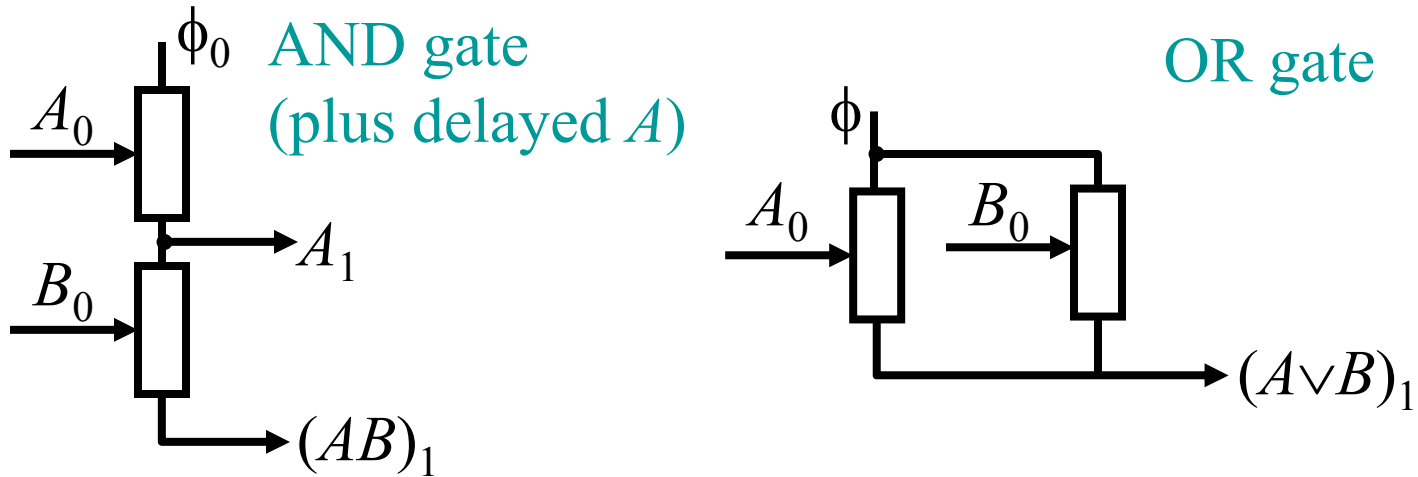




# More Complex Logic Functions



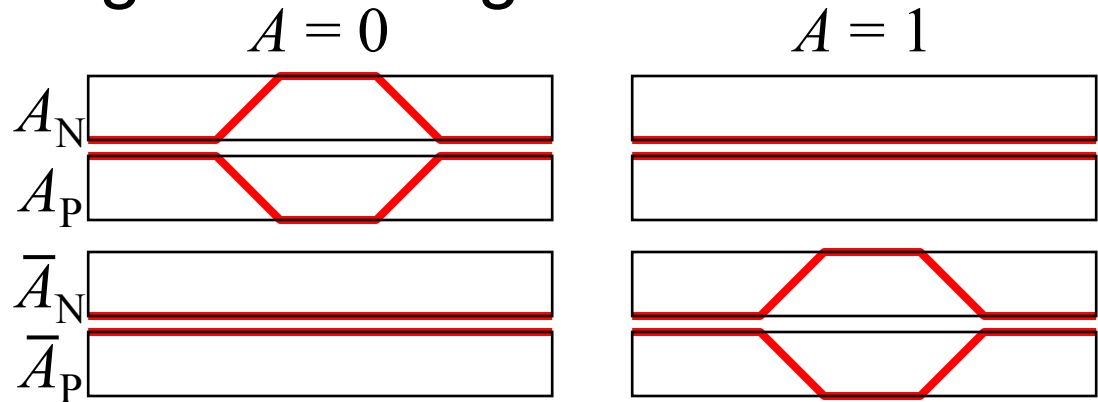
- Non-inverting multi-input Boolean functions:



- One way to do inverting functions in pipelined logic is to use a quad-rail logic encoding:

– To invert, just swap the rails!

- Zero-transistor “inverters.”



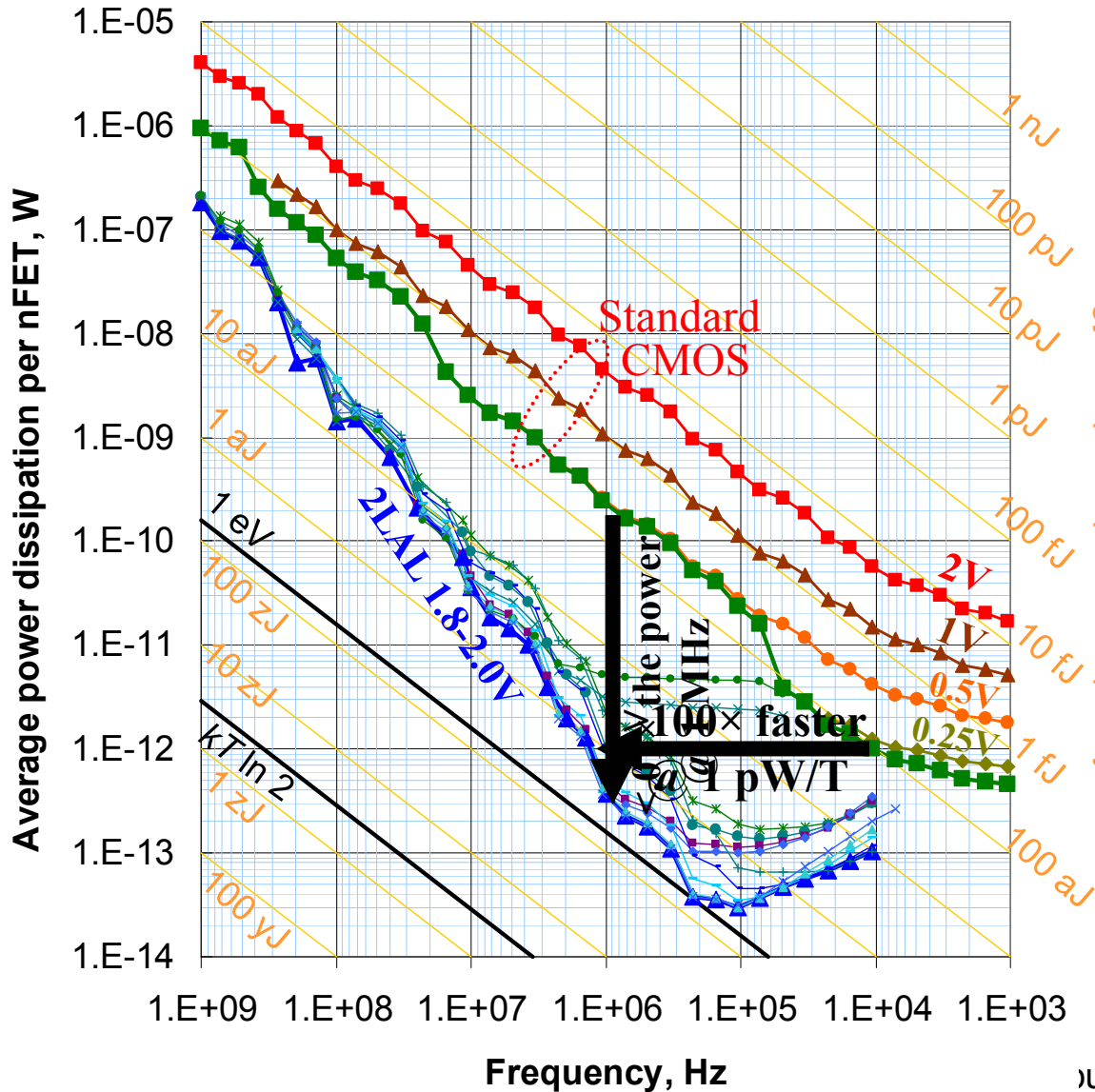




# Simulation Results from Cadence



### Power vs. freq., TSMC 0.18, Std. CMOS vs. 2LAL



### Assumptions & caveats:

- Assumes ideal trapezoidal power/clock waveform.
- Minimum-sized devices,  $2\lambda \times 3\lambda$ 
  - \*  $.18 \mu\text{m (L)} \times .24 \mu\text{m (W)}$
- nFET data is shown
  - \* pFETs data is very similar
- Various body biases tried
  - \* Higher  $V_{th}$  suppresses leakage
- Room temperature operation.
- Interconnect parasitics have not yet been included.
- Activity factor (transitions per device-cycle) is 1 for CMOS, 0.5 for 2LAL in this graph.
- Hardware overhead from fully-adiabatic design style is not yet reflected
  - \*  $\geq 2\times$  transistor-tick hardware overhead in known reversible CMOS design styles

puting"



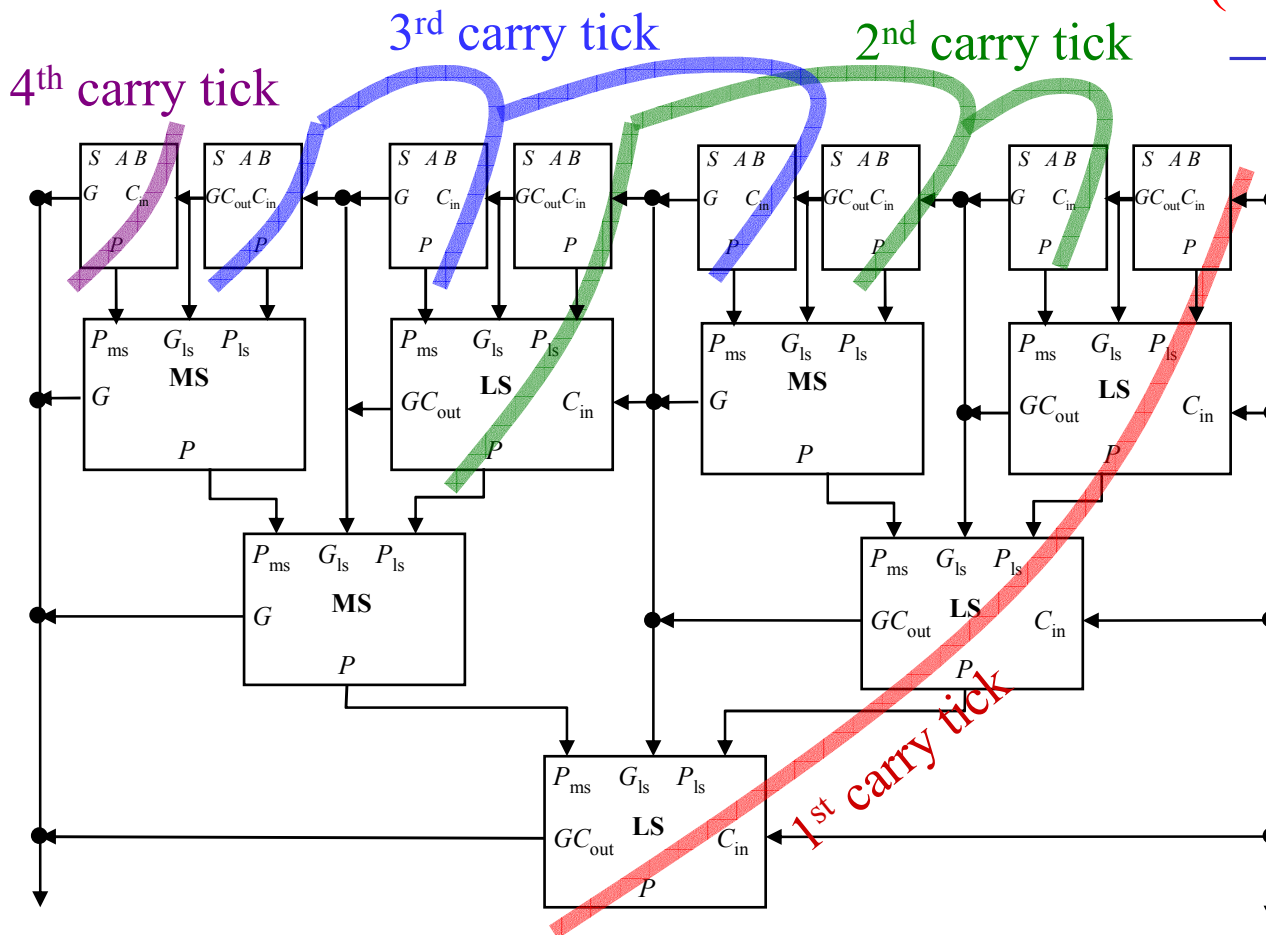
# $O(\log n)$ -time carry-skip adder



(8 bit segment shown)

With this structure, we can do a  $2^n$ -bit add in  $2(n+1)$  logic levels  
→  $4(n+1)$  reversible ticks  
→  $n+1$  clock cycles.

Hardware overhead is  $< 2 \times$  regular ripple-carry.



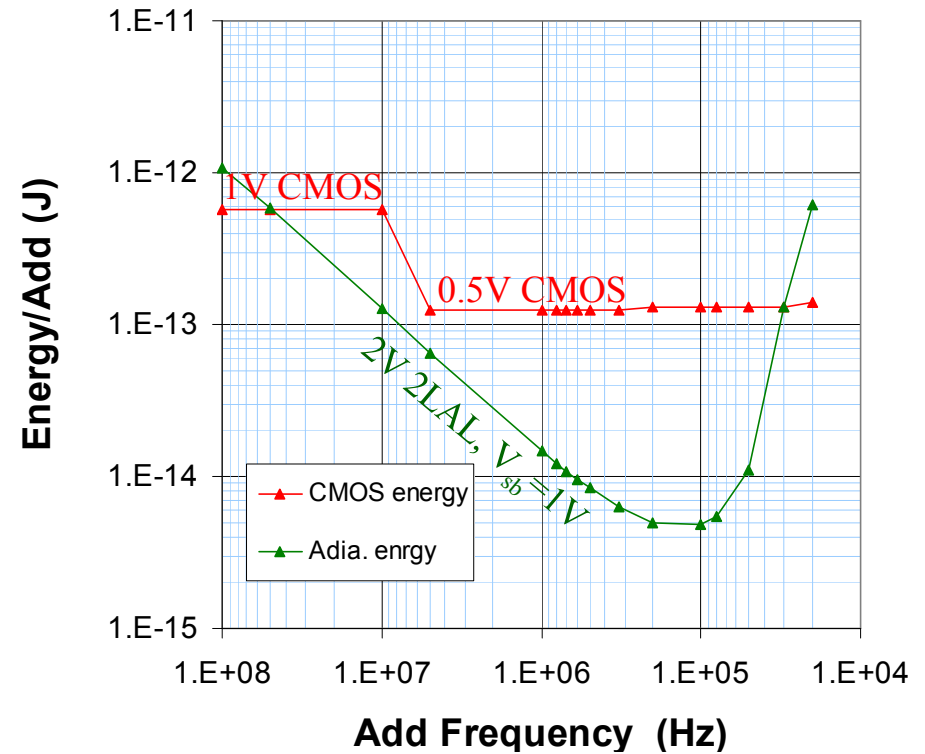
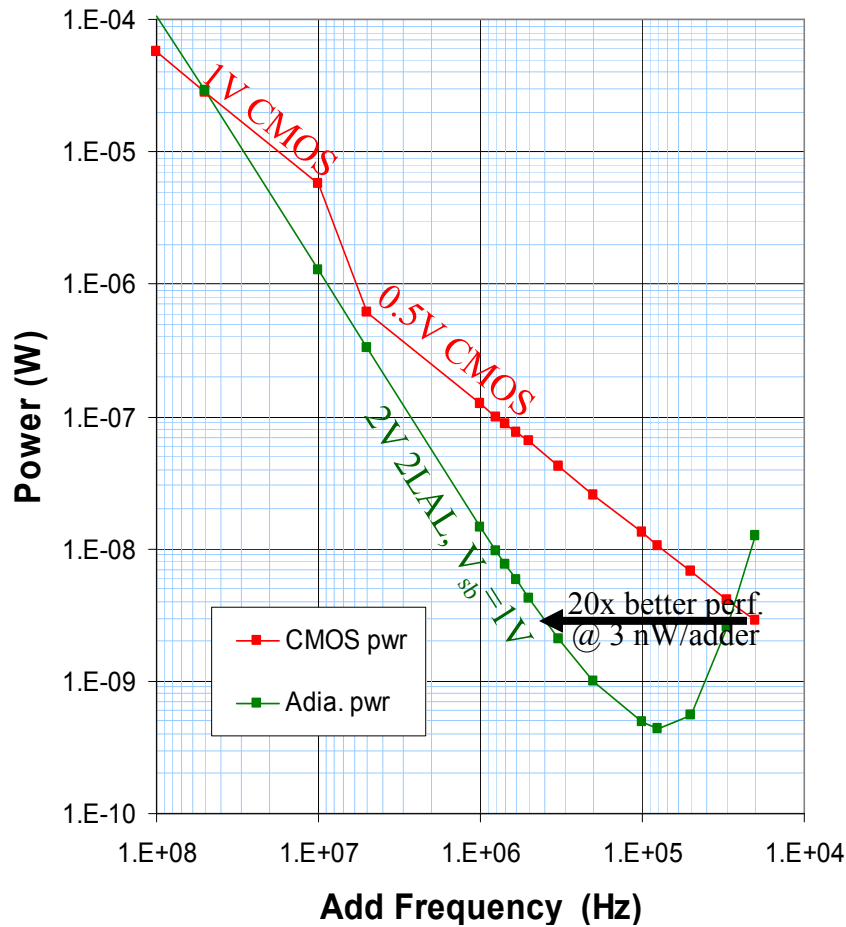


# 32-bit Adder Simulation Results



### 32-bit adder power vs. frequency

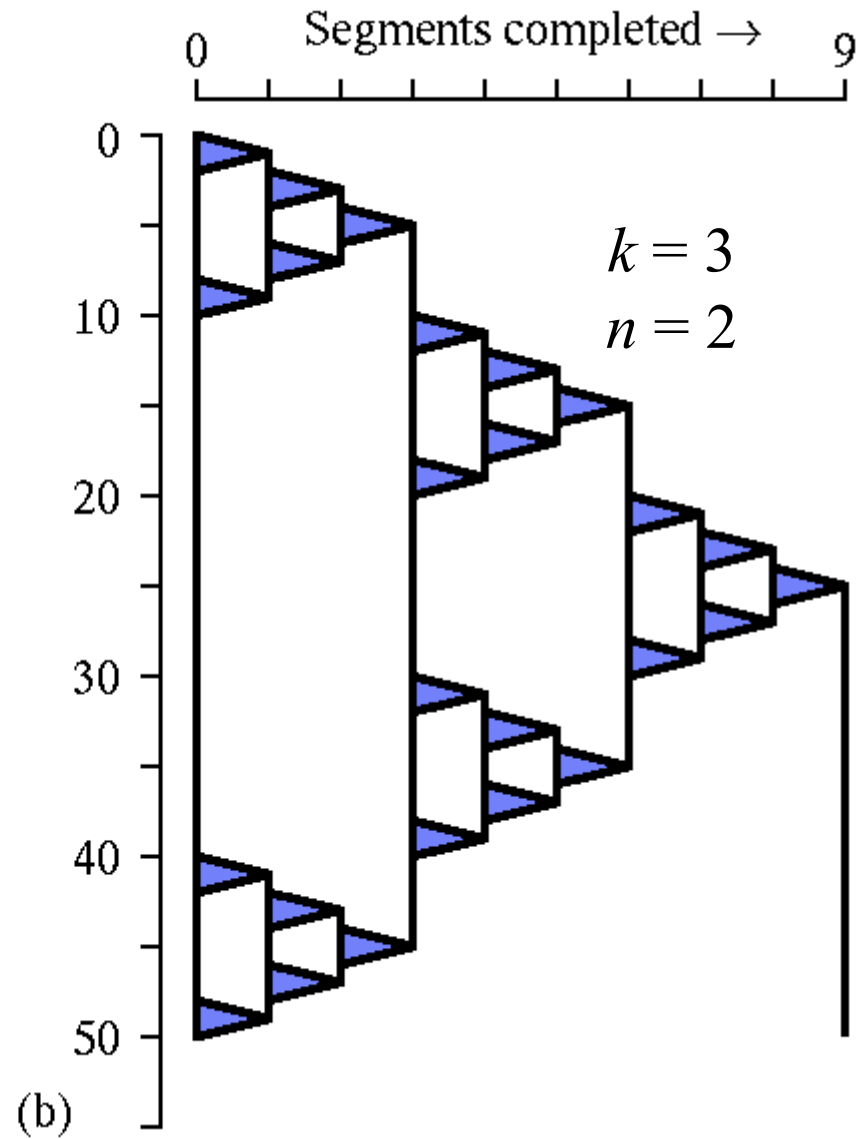
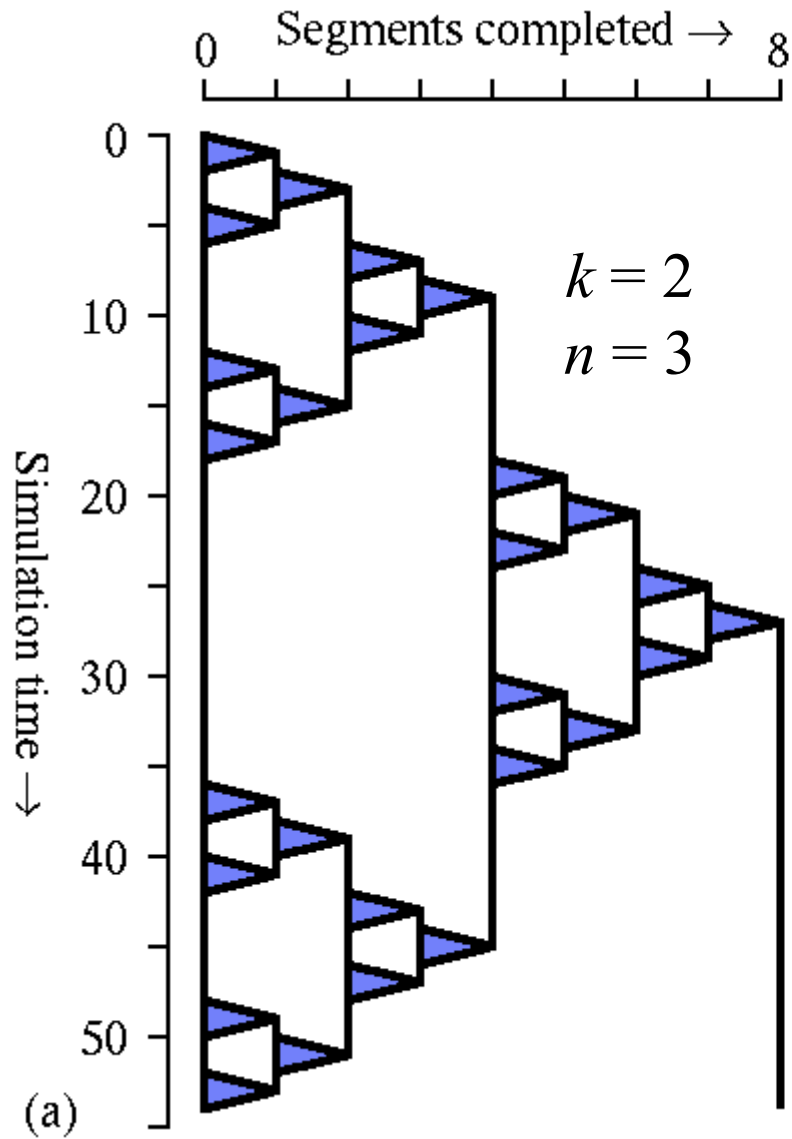
### 32-bit adder energy vs. frequency



(All results normalized to a throughput level of 1 add/cycle)



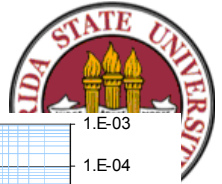
# Bennett '89 algorithm



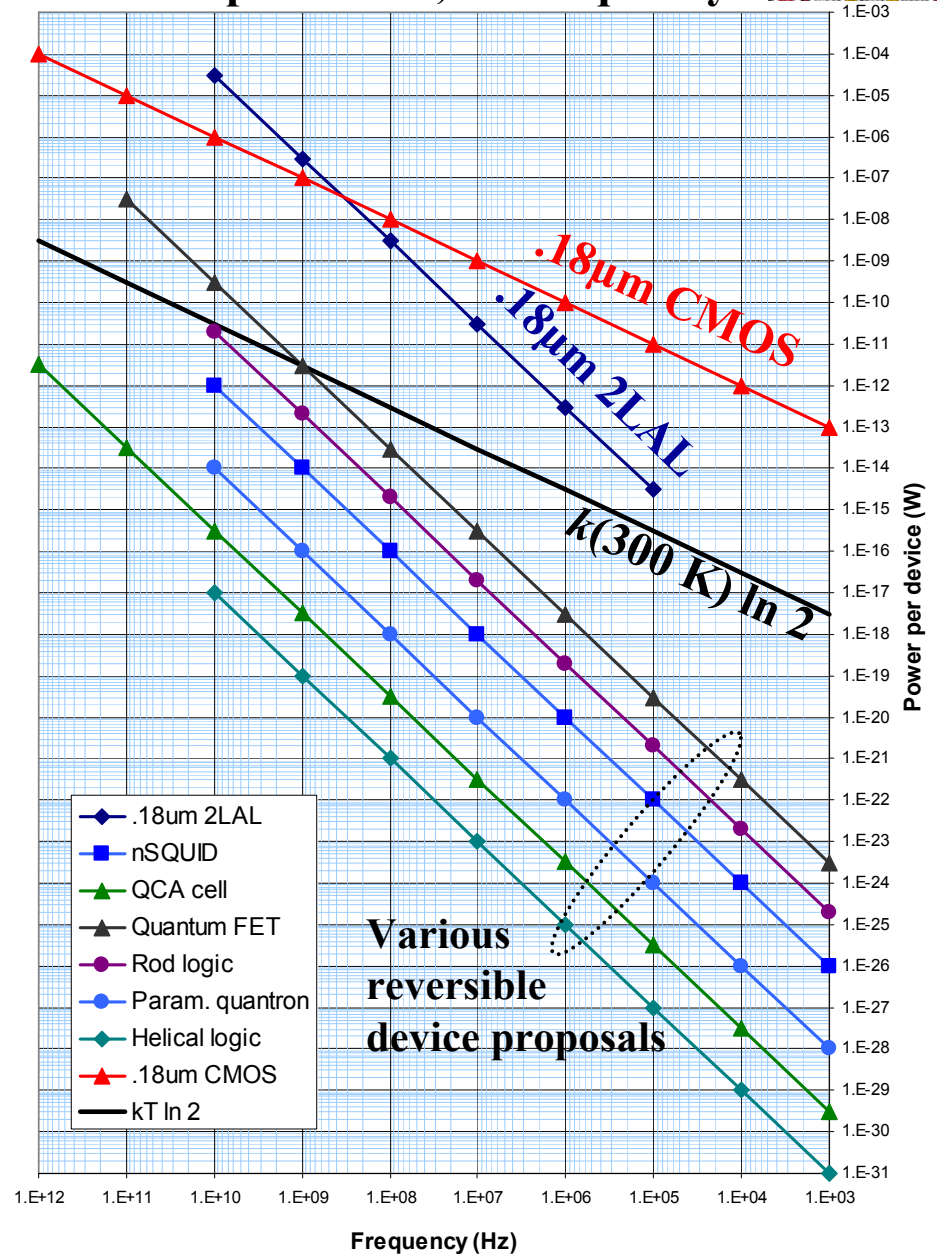


# There's plenty of Room for device improvement...

- Recall, irreversible device technology has at most ~3-4 orders of magnitude of power-performance improvements remaining.
  - And then, the firm  $kT \ln 2$  limit is encountered.
- But, a wide variety of proposed reversible device technologies have been analyzed by physicists.
  - With theoretical power-performance up to 10-12 orders of magnitude better than today's CMOS!
    - Ultimate limits are unclear.



Power per device, vs. frequency





# The Power Supply Problem

- In adiabatics, the factor of reduction in energy dissipated per switching event is limited to (at most) the Q factor of the clock/power supply.

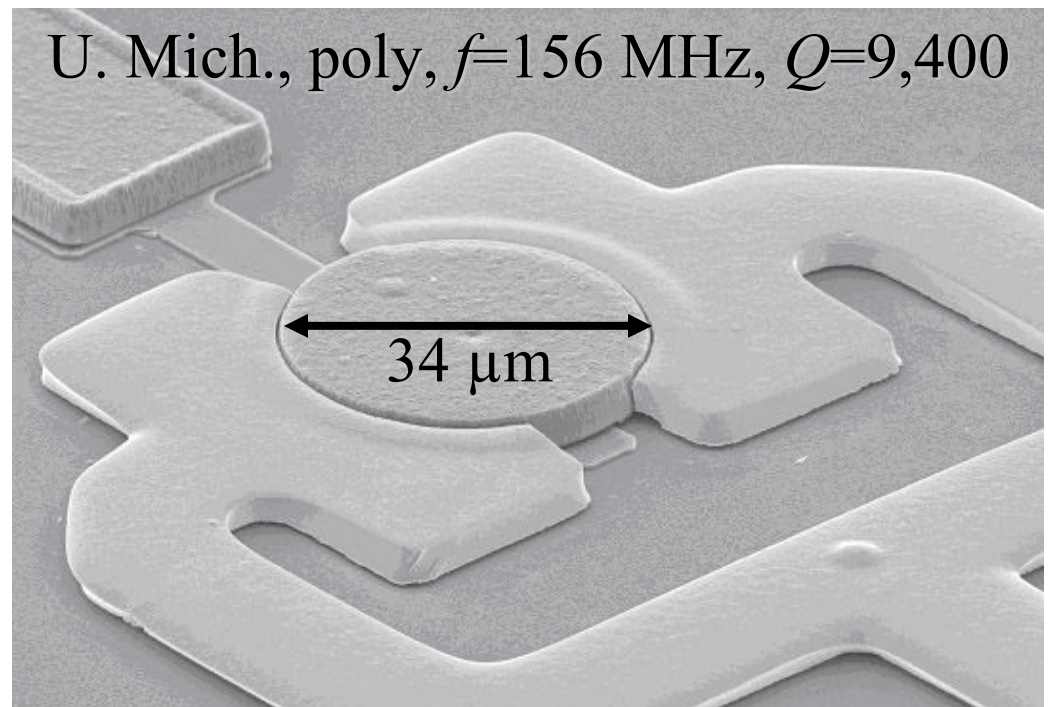
$$Q_{\text{overall}} = (Q_{\text{logic}}^{-1} + Q_{\text{supply}}^{-1})^{-1}$$

- Electronic resonator designs typically have low Q factors, due to considerations such as:
  - Energy overhead of switching a clamping power MOSFET to limit the voltage swing of a sinusoidal LC oscillator.
  - Low coil count and parasitic substrate coupling in typical integrated inductors.
  - Unfavorable scaling of inductor Q with frequency.
- One potential solution that we are presently exploring:
  - Use electromechanical (MEMS) resonators instead!



# MEMS (& NEMS) Resonators

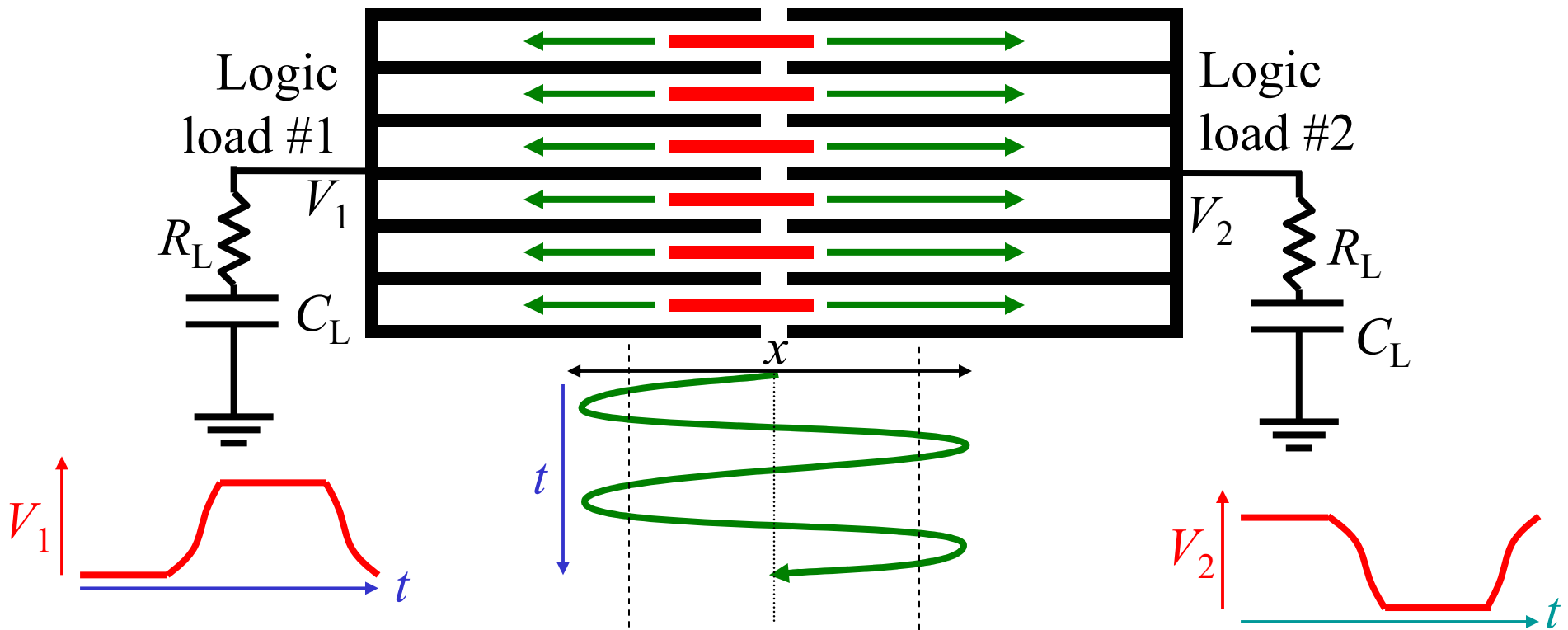
- State of the art of technology demonstrated in lab:
  - Frequencies up to the 100s of MHz, even GHz
  - $Q$ 's  $>10,000$  in vacuum, several thousand even in air!
- An important emerging technology being explored for use in RF filters, *etc.*, in communications SoCs, *e.g.* for cellphones.





# Original Concept

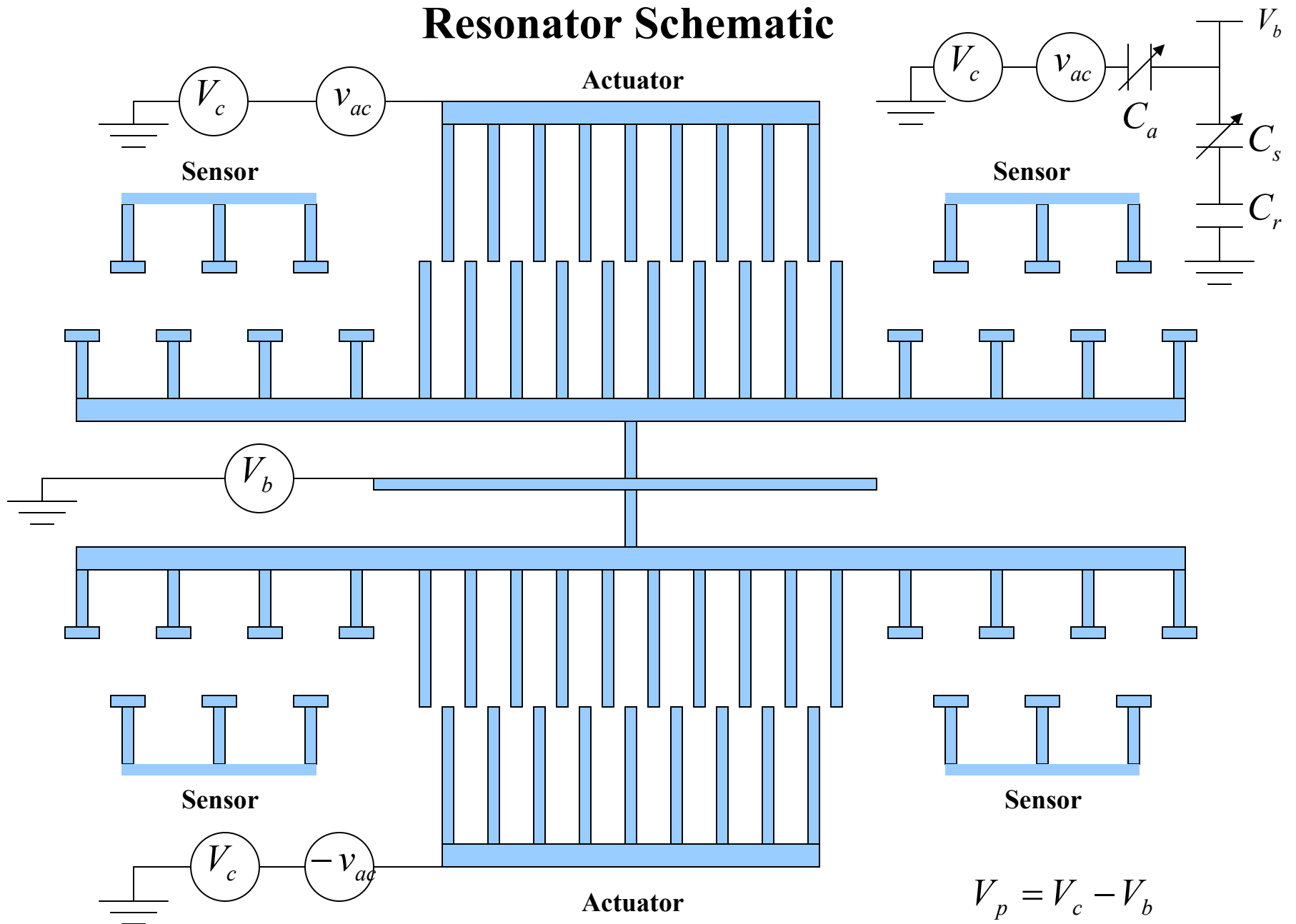
- Imagine a set of charged plates whose horizontal position oscillates between two sets of interdigitated fixed plates.
  - Structure forms a variable capacitor and voltage divider with the load.
- Capacitance changes substantially only when crossing border.
  - Produces nearly flat-topped (quasi-trapezoidal) output waveforms.
  - The two output signals have opposite phases (2 of the 4  $\phi$ 's in 2LAL)





# UF CONFIDENTIAL – PATENT PENDING

## Resonator Schematic





UF CONFIDENTIAL – PATENT PENDING

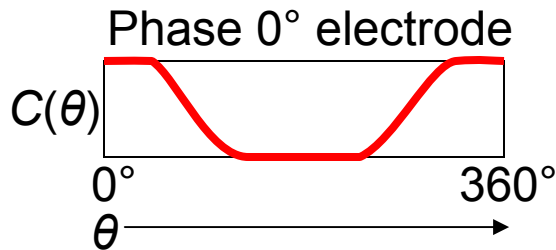
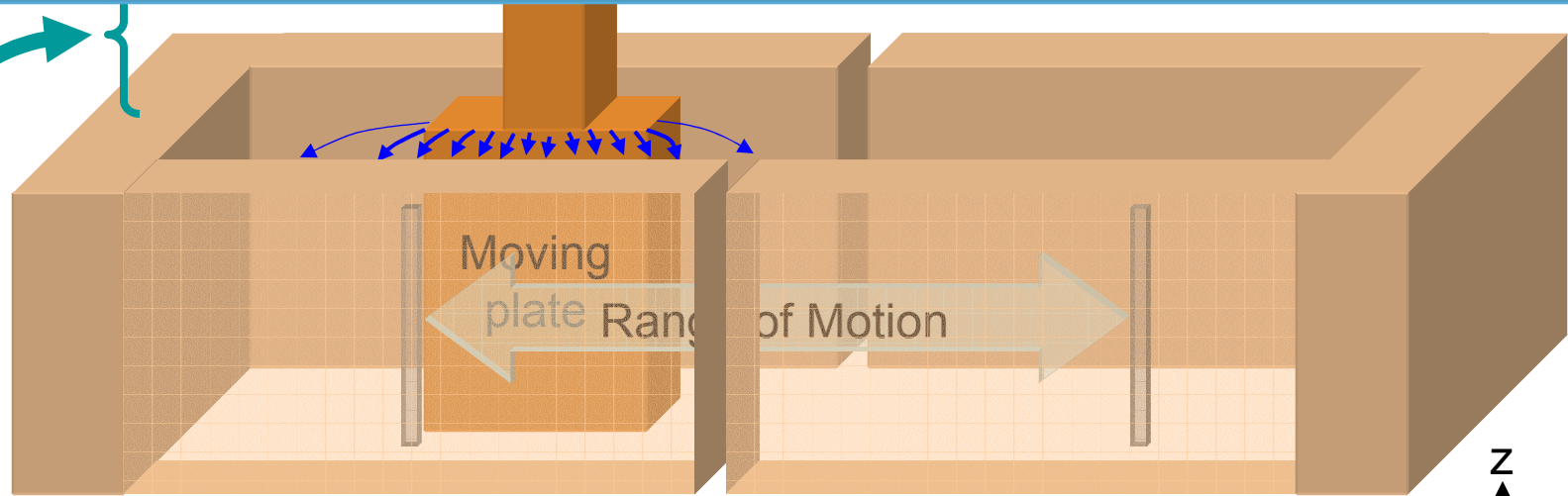


# New Comb Finger Shape IV

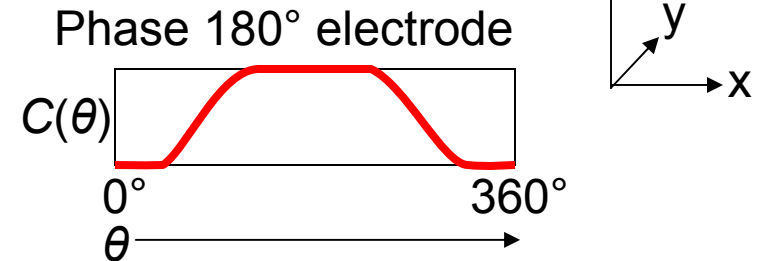
← Arm anchored to nodal points of fixed-fixed beam flexures, located a little ways away, in both directions (for symmetry) →

Moving metal plate support arm/electrode

Is this etch legal?



● Repeat interdigitated structure arbitrarily many times along y axis, all anchored to the same flexure



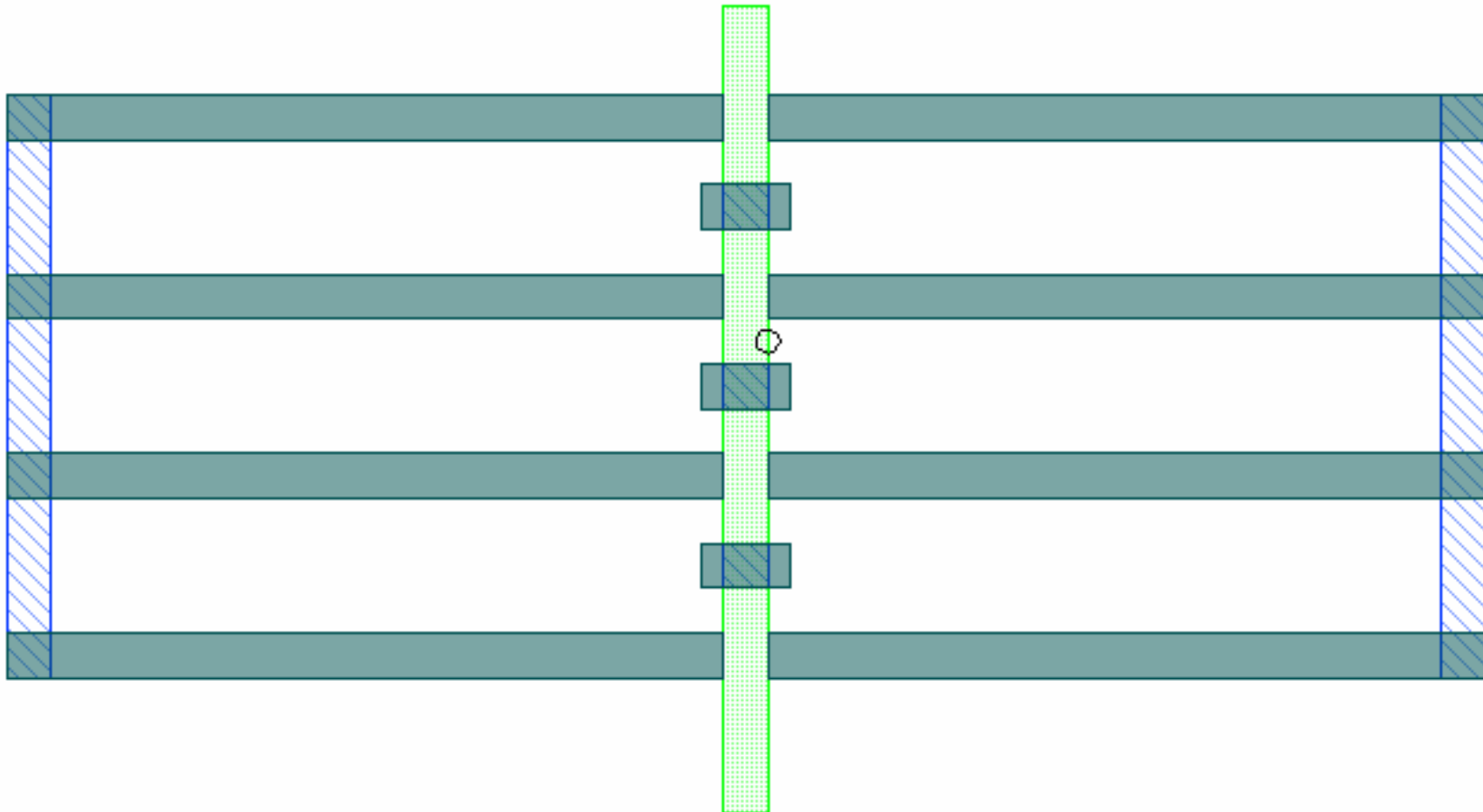
Or, if we can do the structure on the previous slide, then why not this one too? Or, will there be a problem etching the intervening silicon out from in between the metal/oxide layers and the bulk substrate?



**UF CONFIDENTIAL – PATENT PENDING**



# Another Candidate Layout

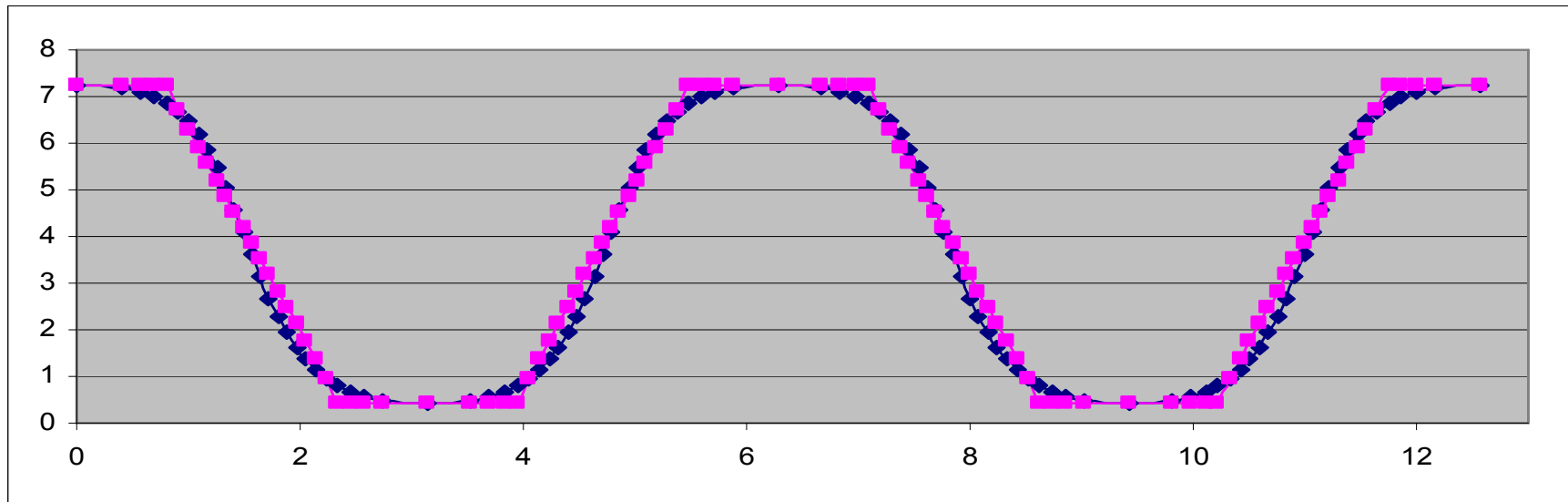
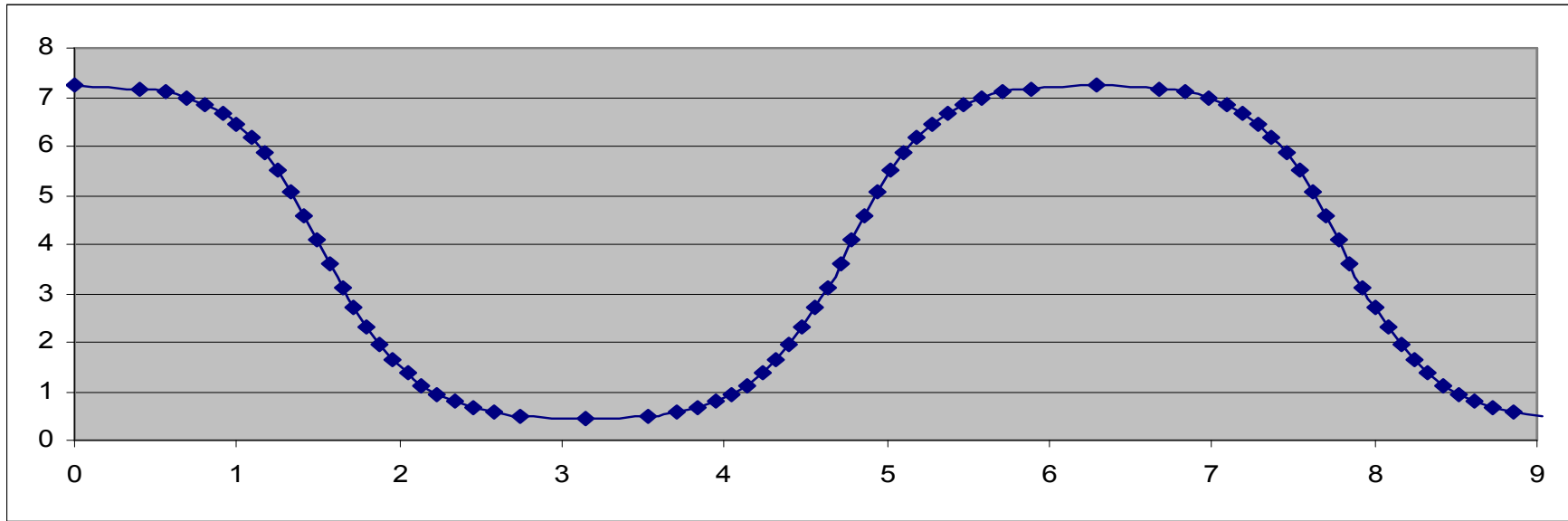




**UF CONFIDENTIAL – PATENT PENDING**

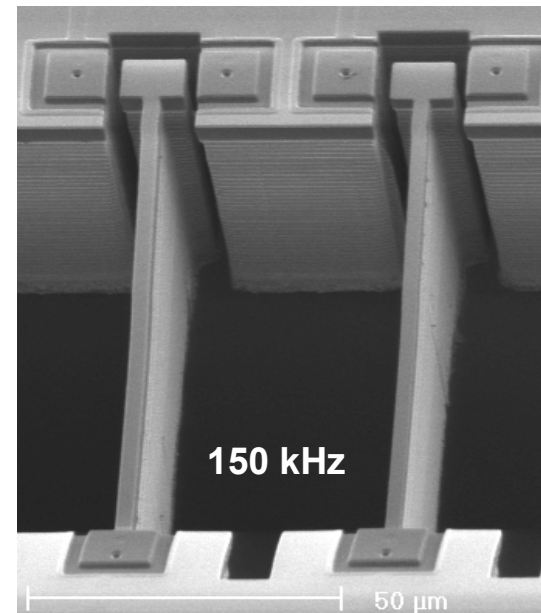
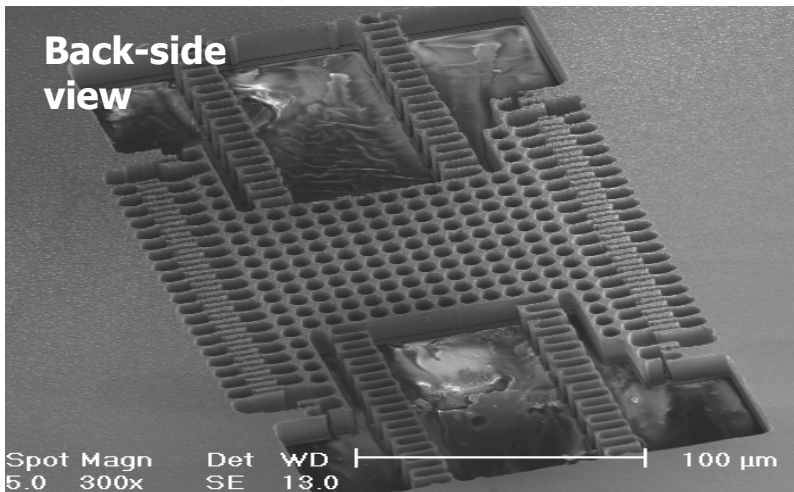
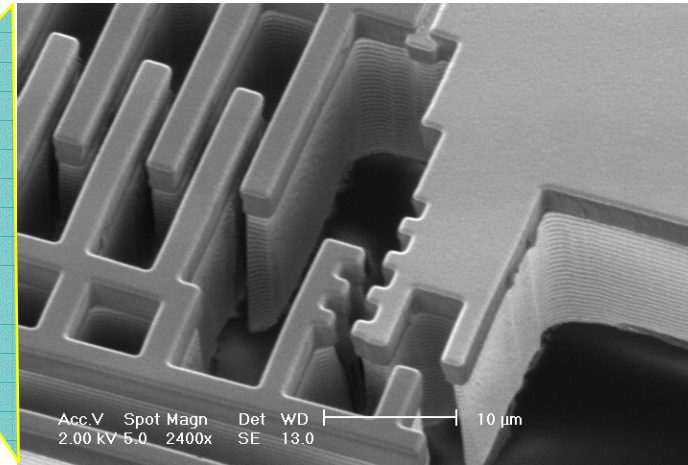
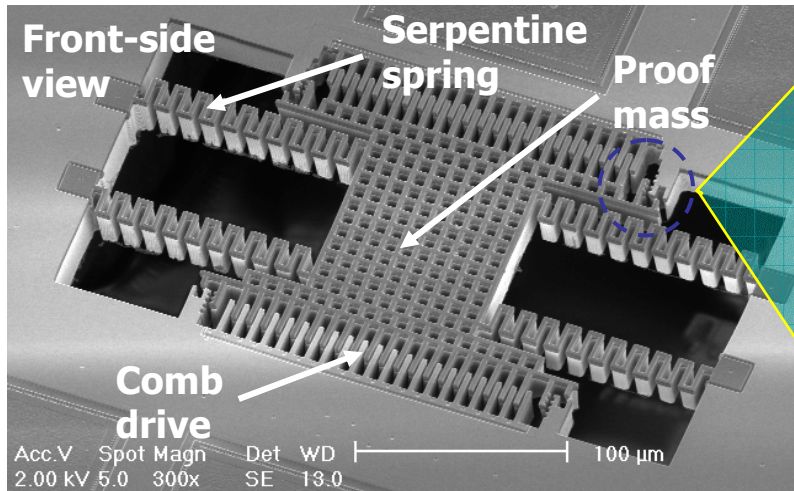


# New simulation results





# DRIE CMOS-MEMS Resonators



## Resonators

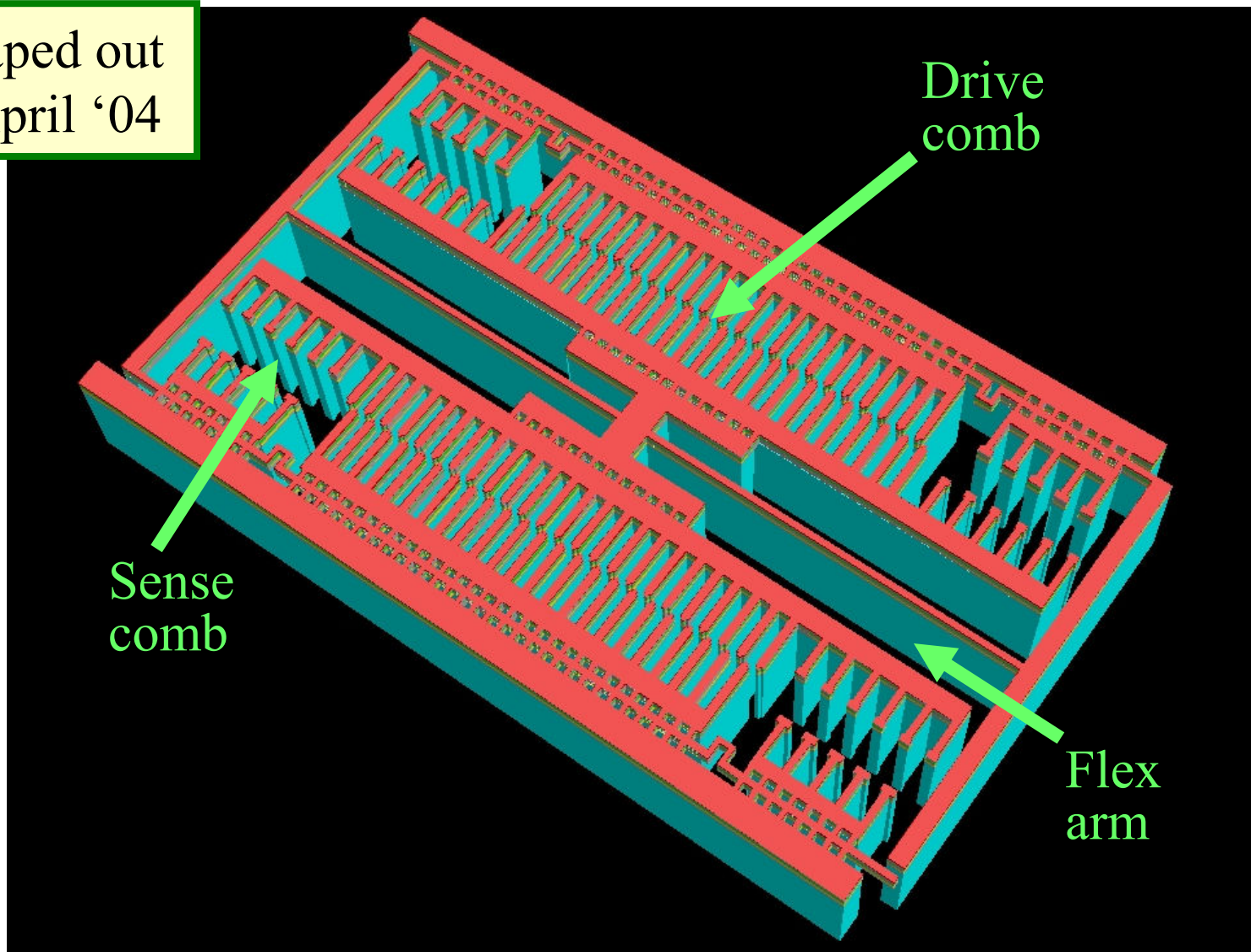


**UF CONFIDENTIAL – PATENT PENDING**



# Post-TSMC35 AdiaMEMS Resonator

Taped out  
April '04



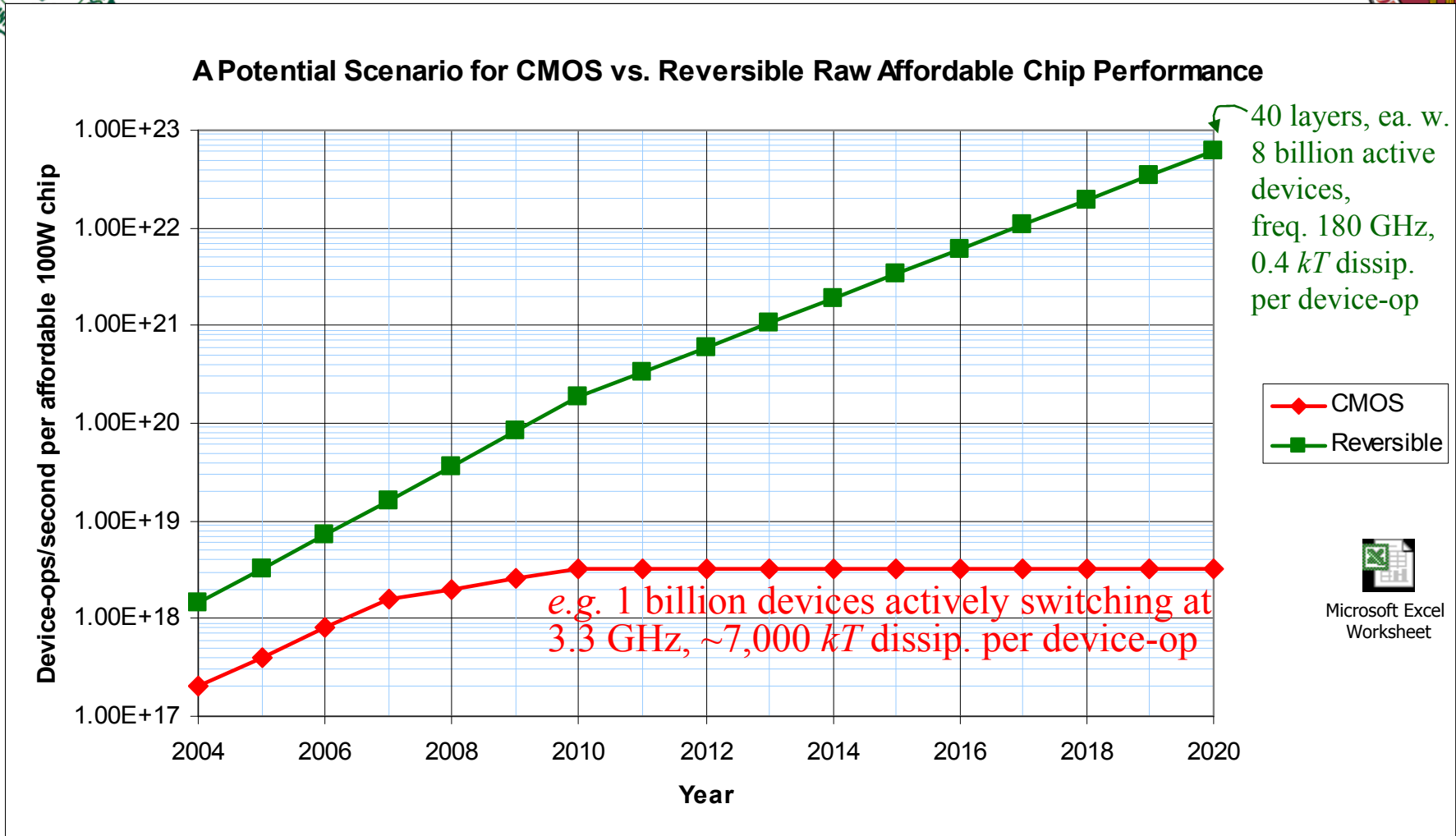


# One Potential Scaling Scenario for Reversible Computing Technology

- Assume *energy coefficient* (energy diss. / freq.) of reversible technology continues declining at historical rate of  $16\times / 3$  years, through 2020.
  - For adiabatic CMOS,  $c_E = CV^2RC = C^2V^2R$ .
    - This has been going as  $\sim \ell^4$  under constant-field scaling.
  - But, requires new devices after CMOS scaling stops.
    - However, many candidates are waiting in the wings...
- Assume number of affordable *layers* of active circuitry per chip (or per package, e.g., stacked dies) doubles every 3 years, through 2020.
  - Competitive pressures will tend to ensure this will happen, esp. if device-size scaling stops, as assumed.



# Result of Scenario



Note that by 2020, there might be as much as a factor of 20,000× difference in *raw* performance per 100W package. (E.g., a 100× overhead factor from reversible design could be absorbed while still showing a 200× boost in performance!)





# Is Reversible Computing Possible?

- This is a worthwhile question to ask, if:
  - By “computing” we mean:
    - scalable, parallel, general-purpose programmable digital computation.
  - By “reversible computing,” we mean:
    - computing with  $\ll E$  energy dissipation per equivalent irreversible logic operation,
      - where  $E$  is the typical minimum logic signal energy
  - And if by “Is it possible?” we mean:
    - Could cost-effective reversible machines be economically manufactured within 20-30 years,
      - Given a sufficient near-term investment in the enabling basic research?



# Status of this Question

- The absolutely most honest scientific answer is:
  - No totally confident, definite answer to this question (yes or no) can be given at present.
- Reversible computing has never been proven to be possible.
  - For that, we would need a validated empirical demonstration of it (on top of a demonstrated manufacturing base), or at least a convincingly very complete and clearly buildable physical model.
    - Demonstrations have been built, but not competitive ones.
    - Physical models have been described, but all are incomplete.
- However, RC has never been proven impossible either.
  - Doing so would require a rigorous proof from consensus physics that somehow addresses all physically possible mechanisms.
    - Various supposed “impossibility” arguments have been offered, but all of them have been riddled with holes and logical fallacies.



# Some Important Next Steps



- Construct a complete quantum mechanical model of a set of high-quality building blocks for reversible computers.
  - Some requirements for these devices:
    - Include a universal set of reversible and irreversible logic ops
    - Extremely low energy coefficient (high Q factor at high frequency)
    - Self-contained (time-independent Hamiltonian, no external drivers)
    - Scalably composable (in 2D and 3D interconnected networks of devices)
    - High reliability (low prob. of soft errors in typical operating environments)
    - Self-synchronizing, at least locally (asynchronous OK between large blocks)
    - Physically realizable Hamiltonian (local, and composable from available physical interactions)
- Run detailed and complete physical simulations of complex digital applications composed of the above building blocks.
  - Validate that unforeseen problems do not arise at higher design levels.
- Show how to implement these building blocks in an economically viable (cost-effective) manufacturing process.
  - Show that the resulting systems would operate in a cost-effective fashion, competitively against conventional designs.
- Migrate supporting tools for new & legacy languages & applications to the new mostly-reversible architectural platforms.



# Conclusion

- Reversible computing is possible...
  - As far as fundamental physics can tell us at the moment.
- It is necessary...
  - To prevent computer performance from stalling within the next 1-3 decades.
- It is technologically challenging...
  - A number of research & engineering problems remain to be solved in order to implement it efficiently...
- We need to aggressively push to solve the remaining problems!
  - In order for reversible computing to be available in time to help us achieve extreme supercomputing within the scope of our careers.