

Prospects for High-Aspect-Ratio FinFETs in Low-Power Logic

Mark Rodwell, Doron Elias

University of California, Santa Barbara

High Aspect Ratio Fins for Low-Power Logic

Fin thickness defined by Atomic layer epitaxy

→ *nm thickness control*

Fin height defined by sidewall growth

→ *200 nm high fins*

Enables ~4 nm fin bodies → 8 nm gate length

10:1 more current per unit die area

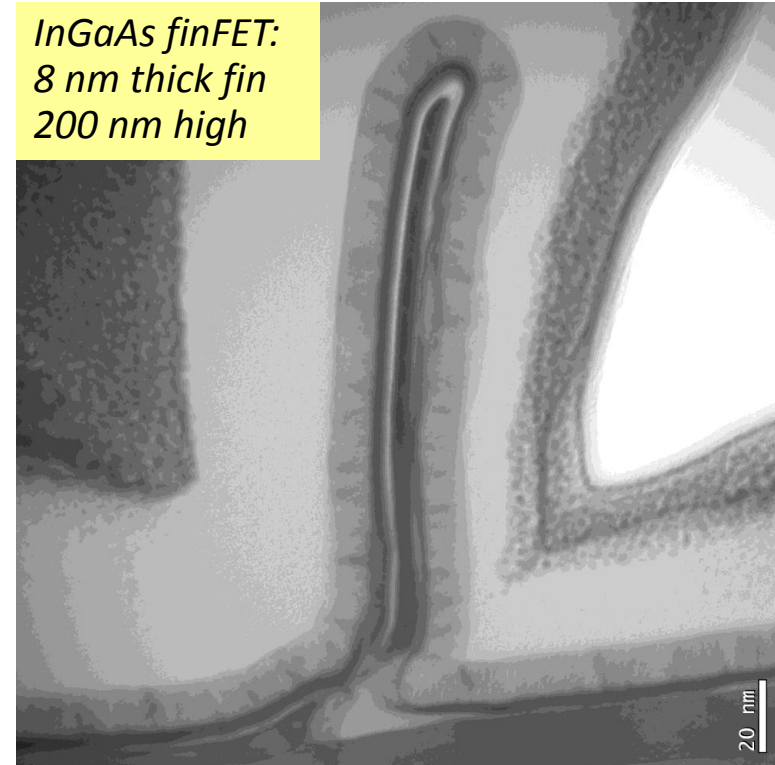
→ *smaller IC die area*

complements lithographic scaling

Enables high speed, ultra low-power logic,

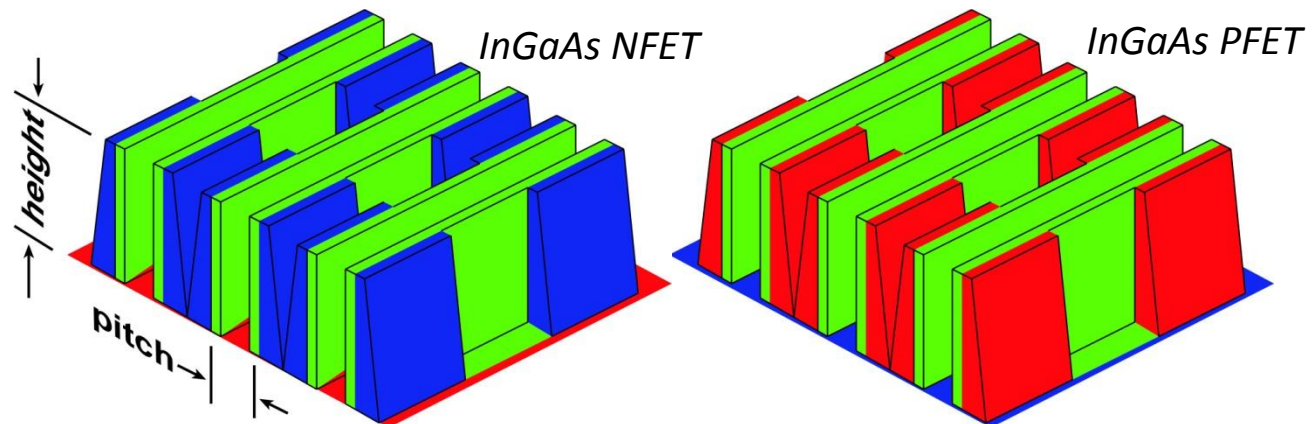
$V_{dd} \sim 300$ mV

*InGaAs finFET:
8 nm thick fin
200 nm high*

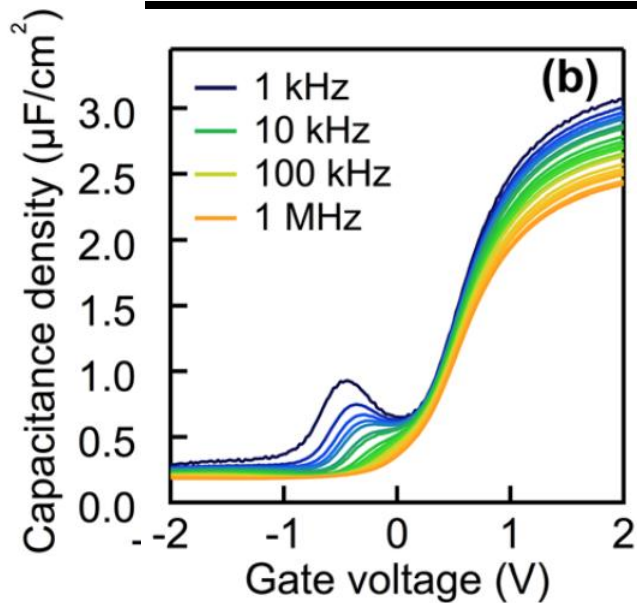


D. Elias, DRC 2013, June, Notre Dame

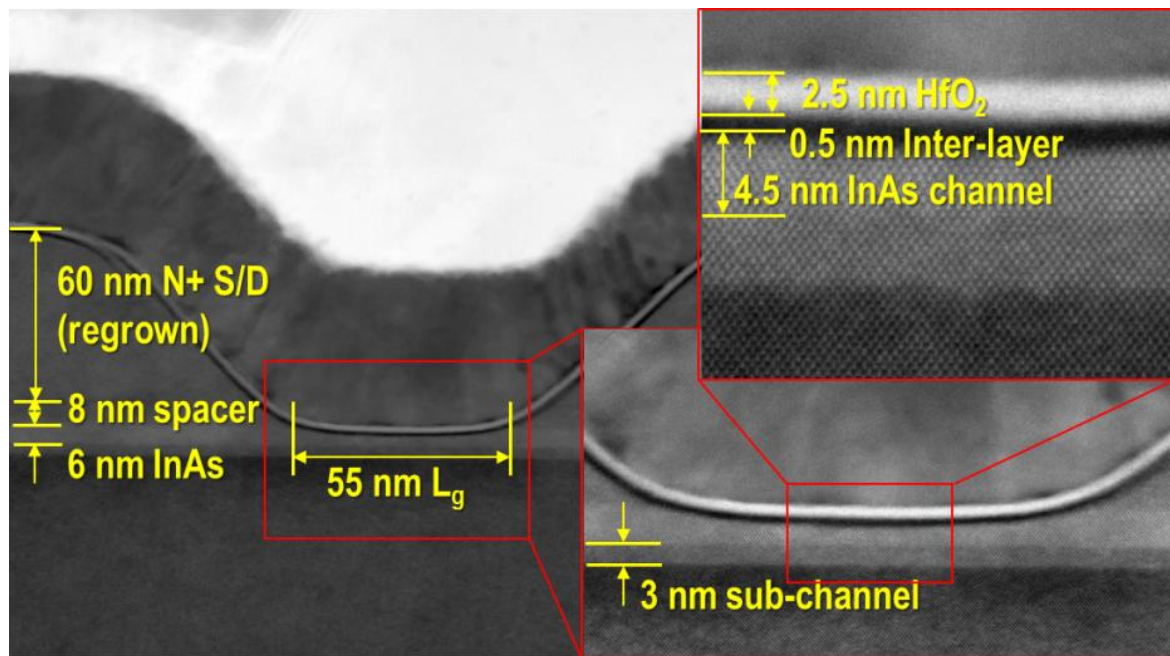
height >> *pitch*



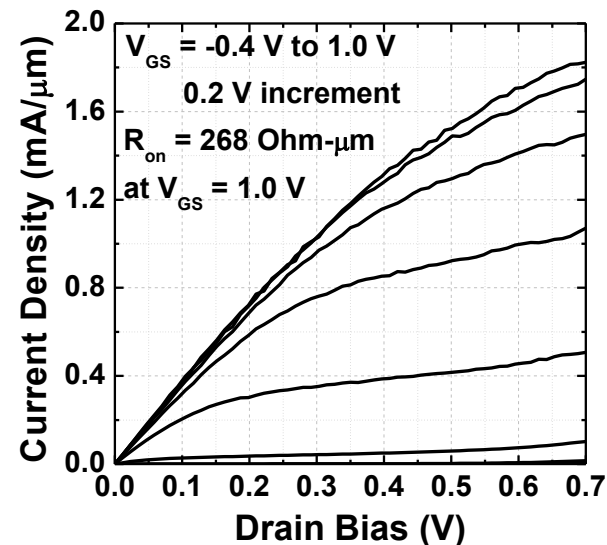
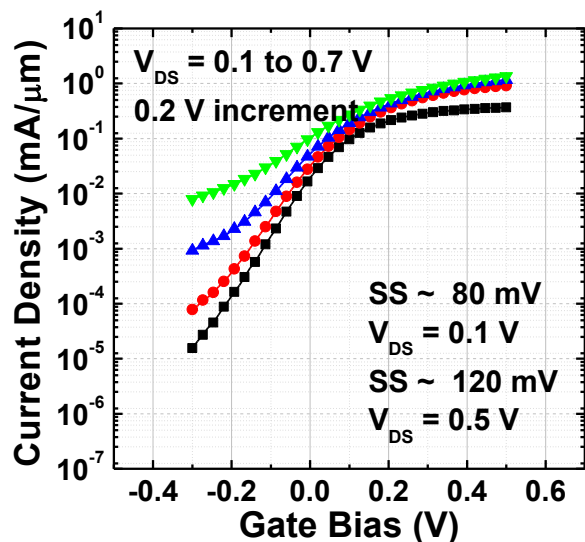
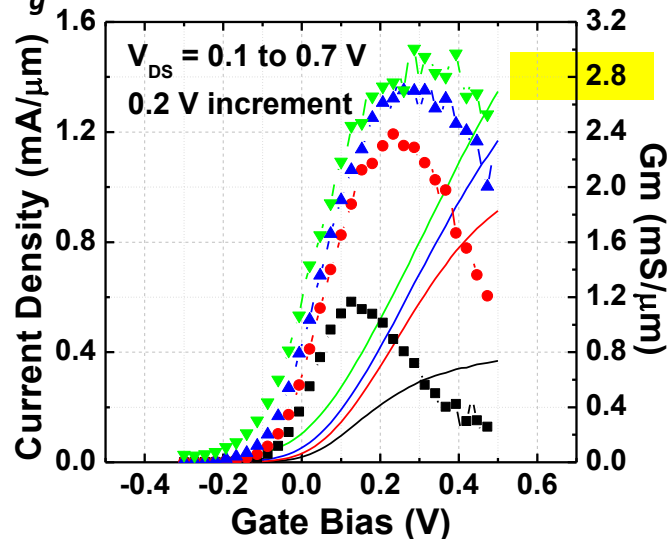
Background: III-V MOS



V. Chobpattana et al (Stemmer group),
APPLIED PHYSICS LETTERS 102, 022907 (2013)



$L_g = 60 \text{ nm}$



FinFETs by Atomic Layer Epitaxy: Why ?

Electrostatics:

body must be thinner than $\sim L_g/2$

→ less than 4 nm thick body for 8 nm L_g

Problem:

threshold becomes sensitive to body thickness

$$\delta V_{th} \propto \delta T_{body} / T_{body}^3$$

Problem:

low mobility unless surfaces are very smooth

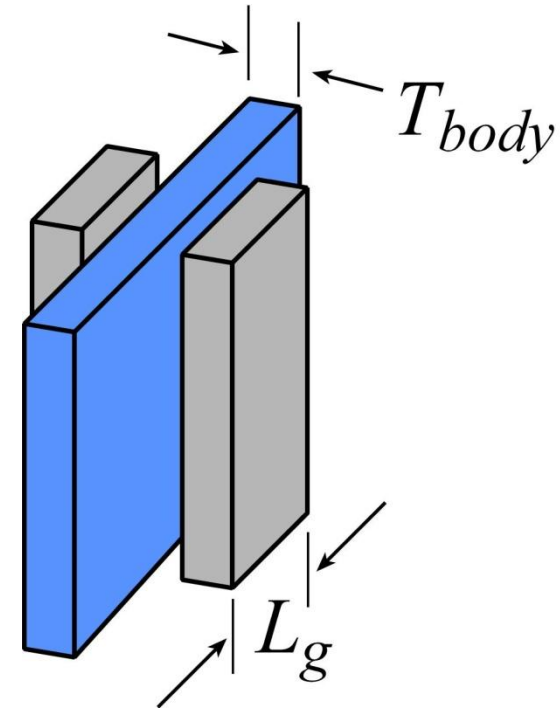
$$\mu \propto T_{body}^6 / \delta T_{body}^2$$

Implication: *At sub-8-nm gate length, need :*

atomically-smooth interfaces

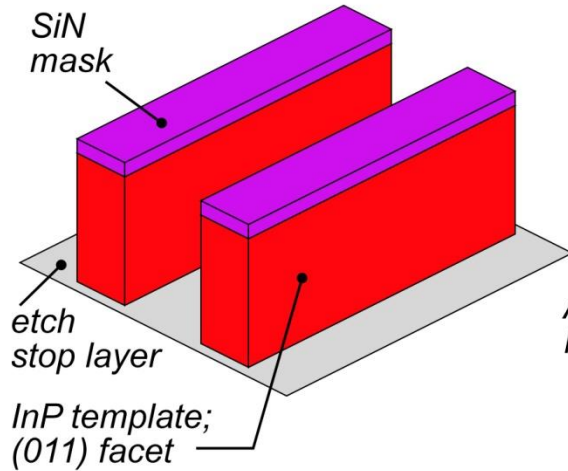
atomically-precise control of channel thickness

side benefit: high drive current → low-voltage, low-power logic

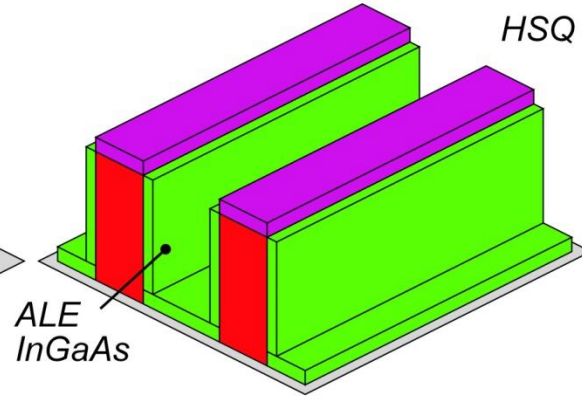


ALE-Defined finFET: Process Flow

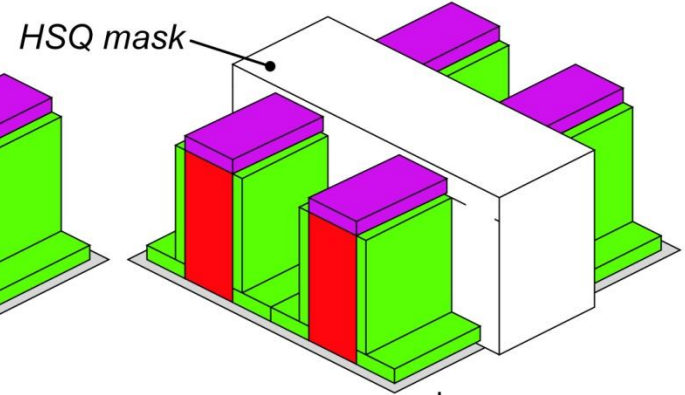
fin template



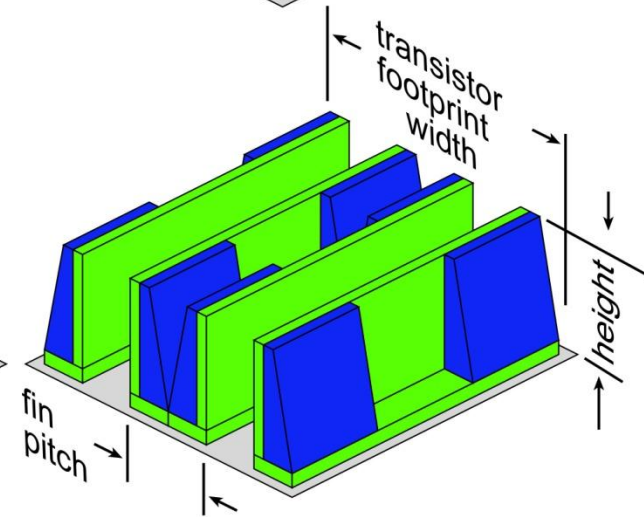
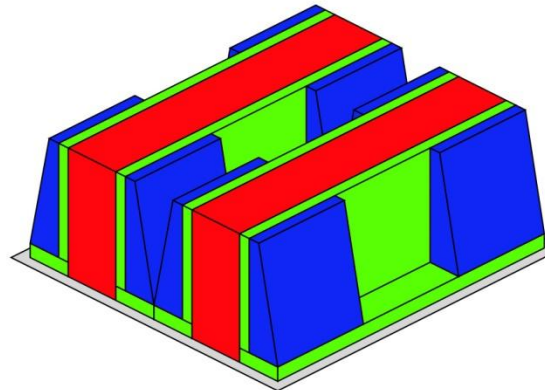
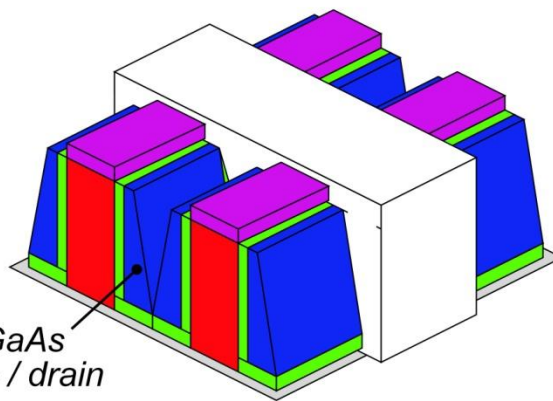
channel ALE



dummy gate



N+ InGaAs source / drain



S/D regrowth

remove masks

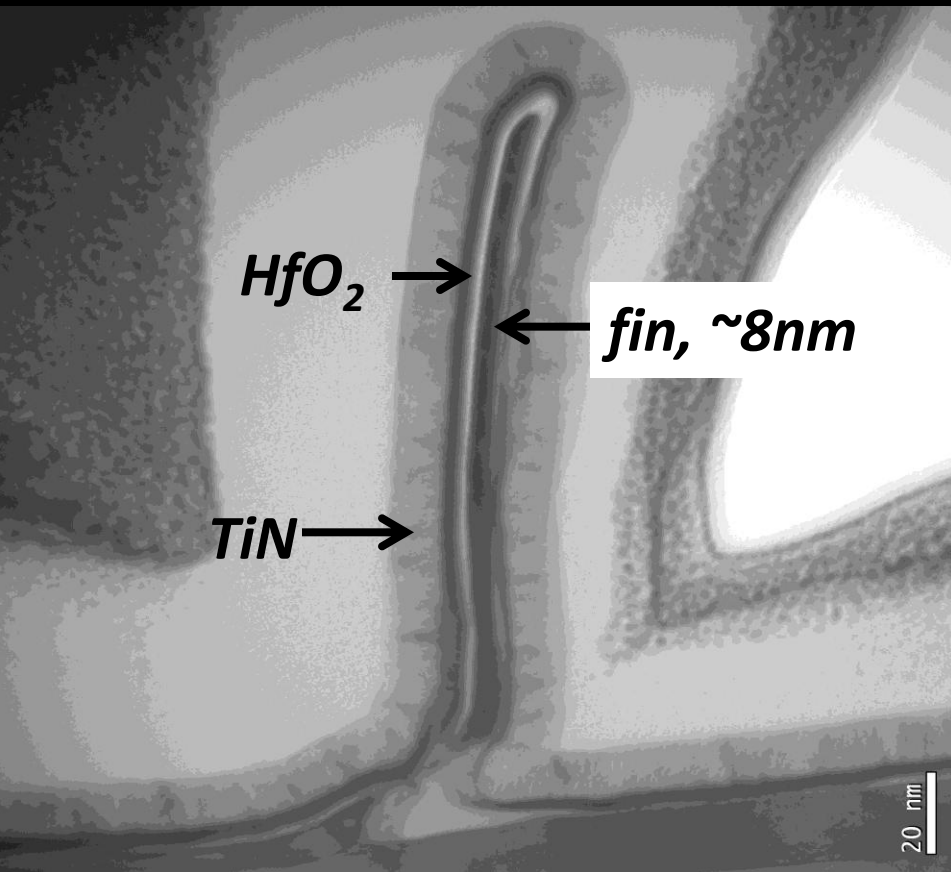
release fins

Fin template: formed by {110}-facet-selective etch → atomically smooth

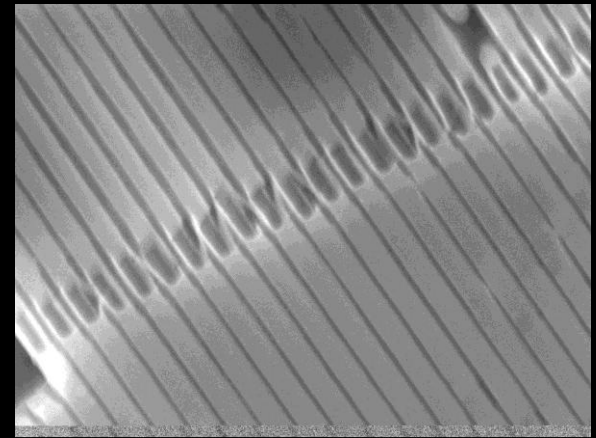
Channel thickness set by ALE growth → atomically precise

Not shown: gate dielectric, gate metal, S/D metal

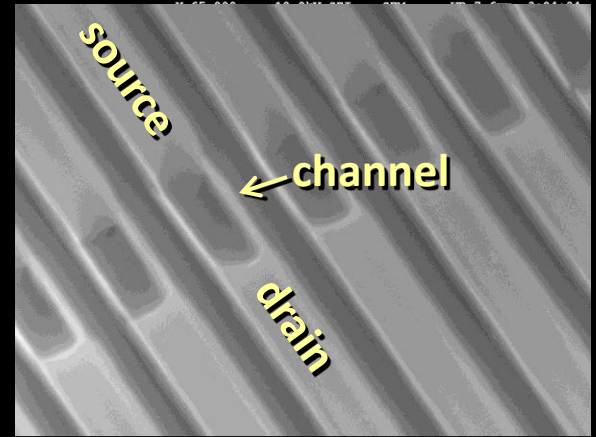
Images



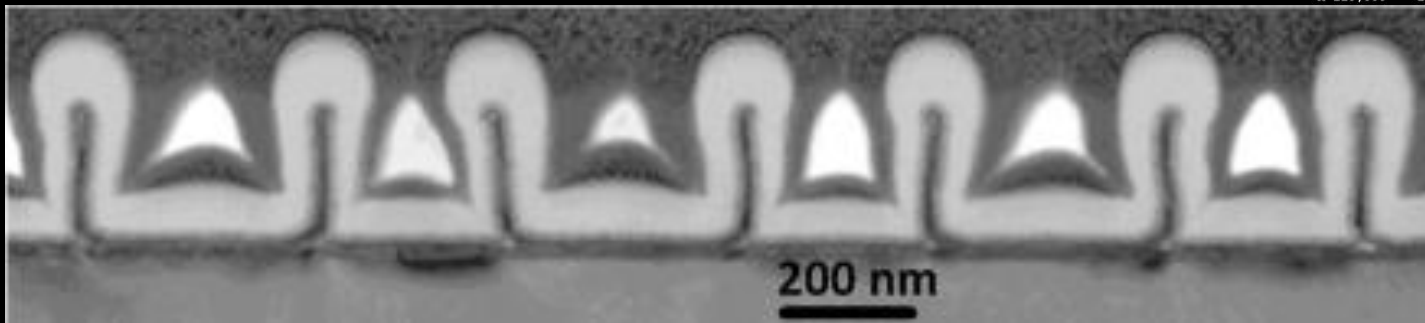
50 nm fin pitch



100 nm fin pitch

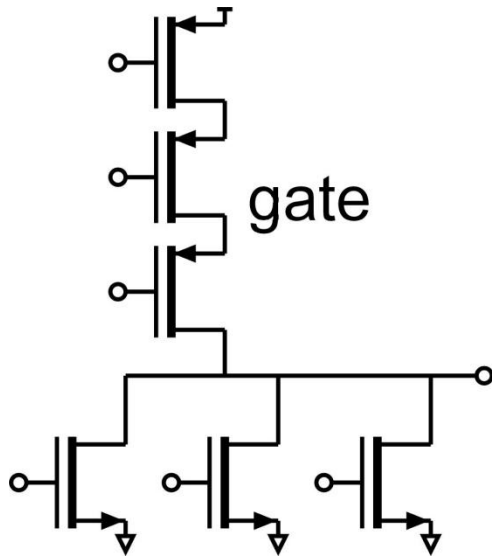
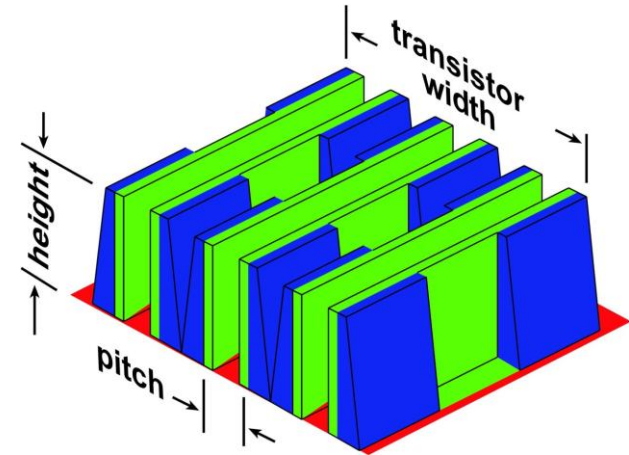


10 nm thick fins, 100 nm tall

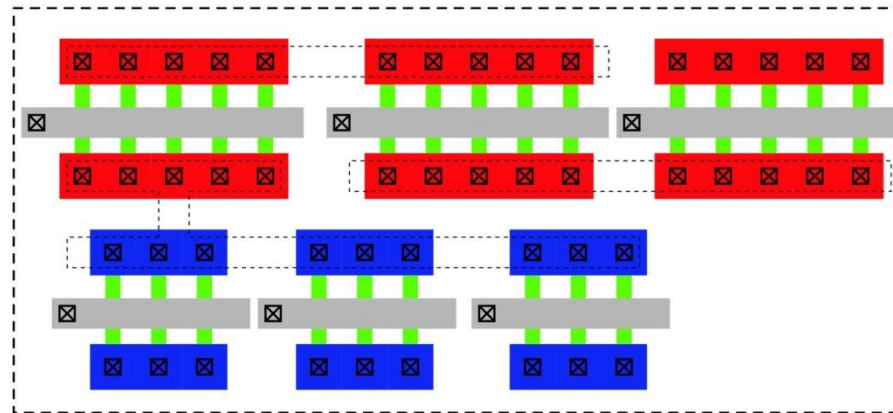


Goal: Tall Fins for High Drive Current

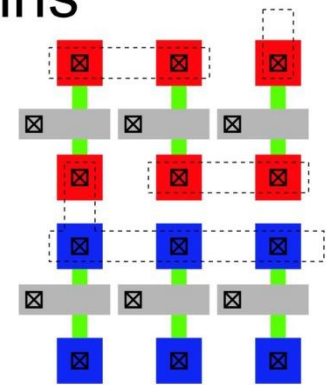
$$\frac{\text{current}}{\text{transistor width}} = J_{\text{surface}} \cdot \frac{\text{fin height}}{\text{fin pitch}}$$



low-current fins



high-current fins

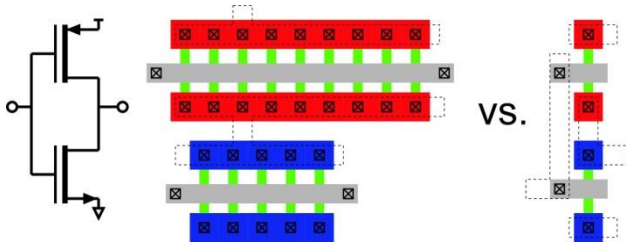


**Goal: fin height \gg fin pitch (spacing) \rightarrow more current per fin
 \rightarrow less fins needed \rightarrow higher integration density**

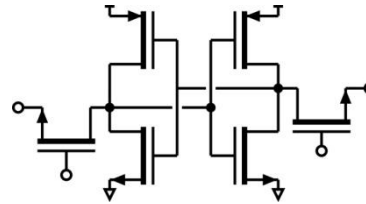
Higher density \rightarrow shorter wires \rightarrow smaller $C_{\text{wire}} V_{\text{dd}}/I$, $C_{\text{wire}} V_{\text{dd}}^2/2$

Is the IC Area Reduction Significant ?

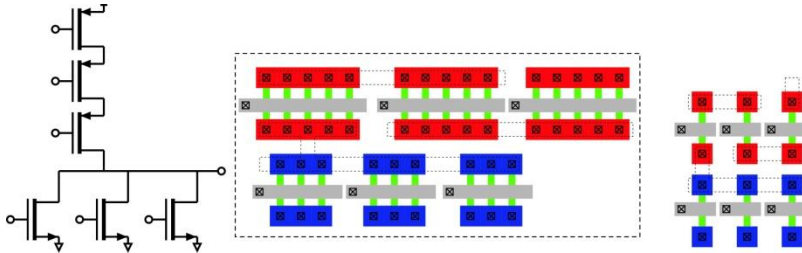
Clock/interconnect drivers need large drive currents.
Area reduction for these is likely substantial.



FETs in **Cache Memory & Registers** are drawn at minimum width
No area reduction for these.



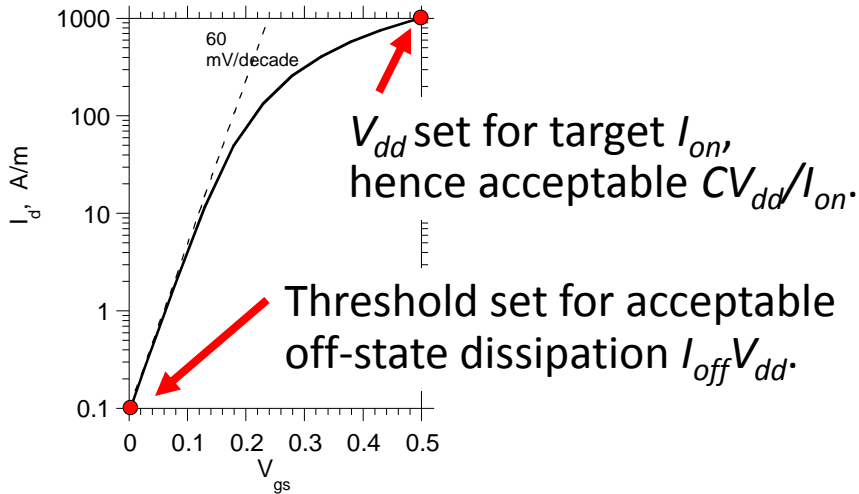
Most, but not all, **Logic Gates** will be drawn at minimum width.



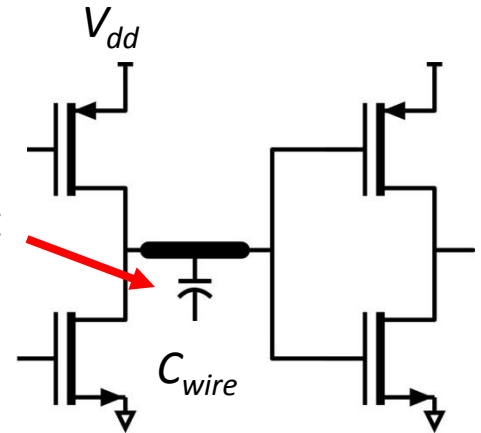
Benefit must be evaluated by VLSI architect, not by device physicist.

300 mV Logic: Can We Address The $CV^2/2$ Limit ?

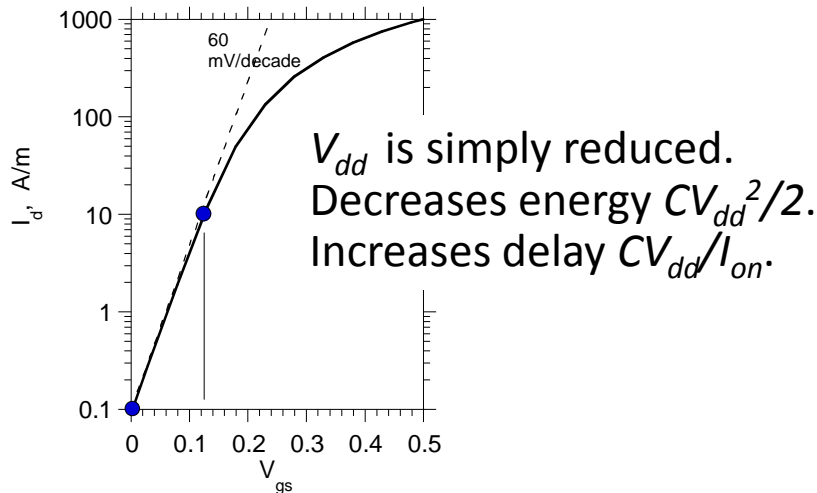
The $CV^2/2$ dissipation limit



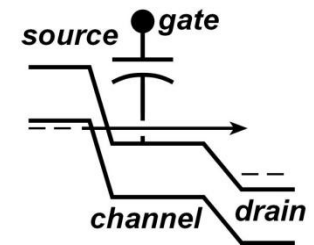
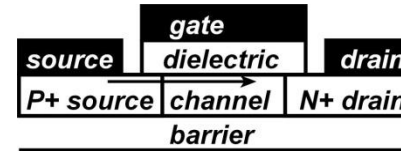
With minimum C_{wire} , a minimum switching energy $C_{wire}V_{dd}^2/2$ is set



Subthreshold logic



Tunnel FETs

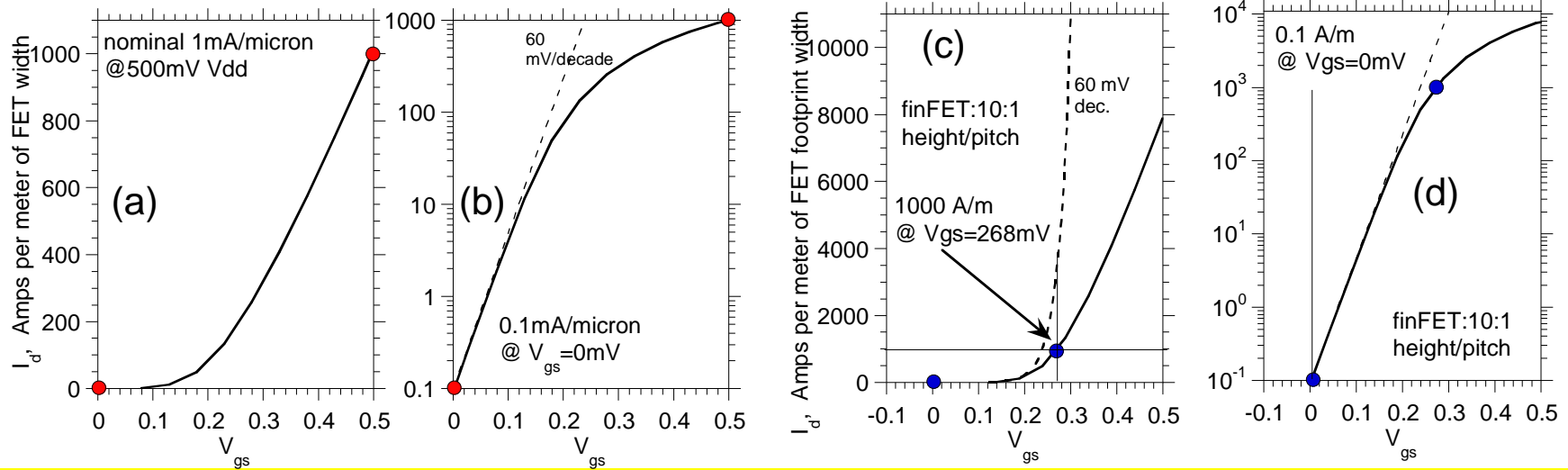


Bandgap of P+ source truncates thermal distribution.

Potential for low I_{off} at low V_{dd} .

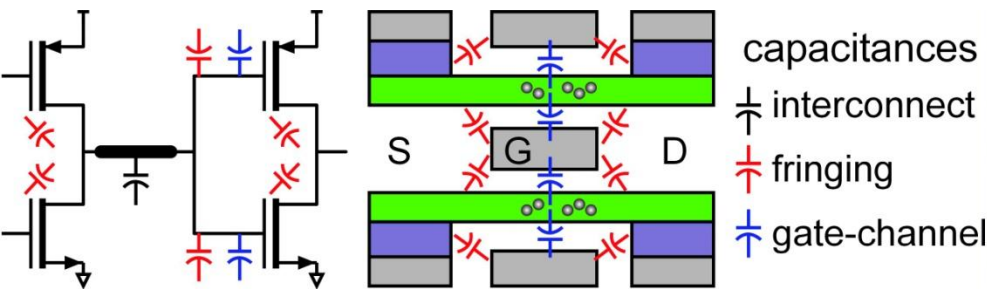
Obtaining high I_{on}/V_{dd} is the challenge.

Goal: Tall Fins for Low-Power, Low-Voltage Logic



Supply reduced from 500mV to 268 mV while maintaining high speed.

3.5:1 power savings ? Must consider FET capacitances.



$$C_{total} \approx 0.2 \text{ fF}/\mu\text{m} \cdot L_{wire} \quad \text{interconnect}$$

$$+ 15 \cdot 0.3 \text{ fF}/\mu\text{m} \cdot W_g \quad \text{fringing}$$

$$+ 3 \cdot I_{on} L_g / v_{injection} V_{DD} \quad \text{gate - channel}$$

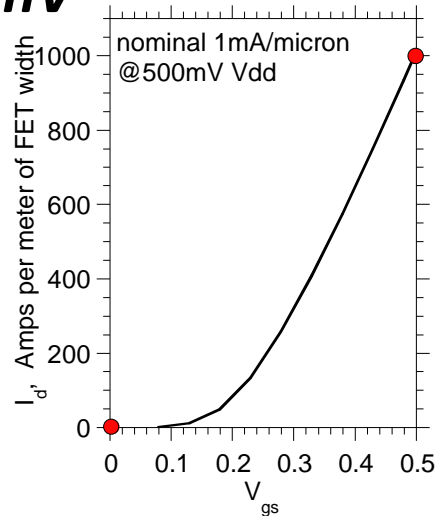
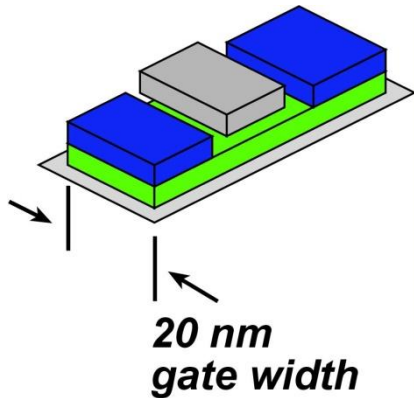
$$\tau_{gate} \approx C_{total} V_{DD} / 2I_{on} \quad \text{delay}$$

$$E_{sw} \approx C_{total} V_{DD}^2 / 2 \quad \text{energy}$$

Assumes (Hodges & Jackson, 2003): (1) Charge-control analysis (2) $I_{on,PFET} / W_g = 0.5 \cdot I_{on,NFET} / W_g$ (3) FO=FI=1

Power and Delay Comparison

Planar FET, $V_{dd}=500\text{ mV}$



$$I_{on}=20\ \mu\text{A}, I_{off}=2\text{nA}$$

$$C_{g-ch}=I_{on}L_g/v_{inj}V_{dd}=1.3\ \text{aF}$$

$$C_{gd-f}=C_{gs-f}=6\ \text{aF}$$

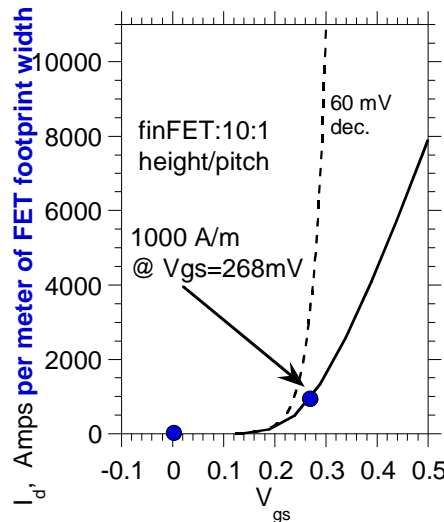
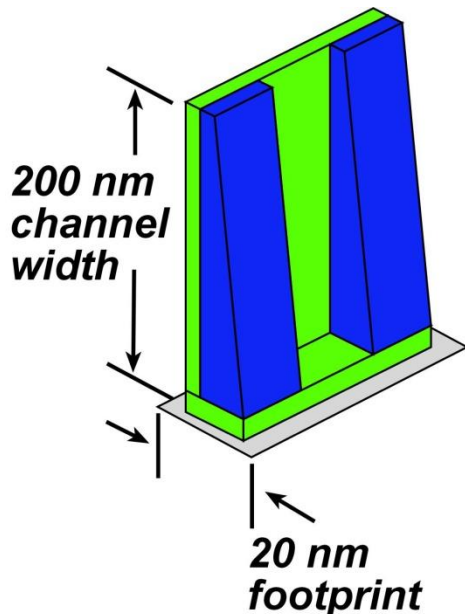
$$C_{wire}=2\ \text{fF}\ (10\ \mu\text{m length})$$

$$C_{total}=2.1\ \text{fF}\ (\text{various multipliers})$$

$$\text{delay}=52\ \text{ps}$$

$$C_{total}V_{DD}^2=0.26\ \text{fJ}$$

tall finFET, $V_{dd}=268\text{ mV}$



$$I_{on}=20\ \mu\text{A}, I_{off}=2\text{nA}$$

$$C_{g-ch}=I_{on}L_g/v_{inj}V_{dd}=3.7\ \text{aF}$$

$$C_{gd-f}=C_{gs-f}=60\ \text{aF}$$

$$C_{wire}=2\ \text{fF}\ (10\ \mu\text{m length})$$

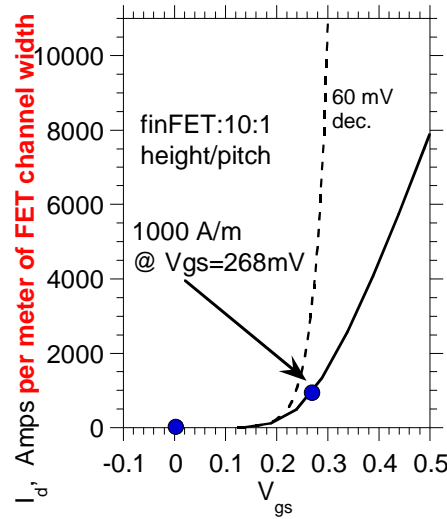
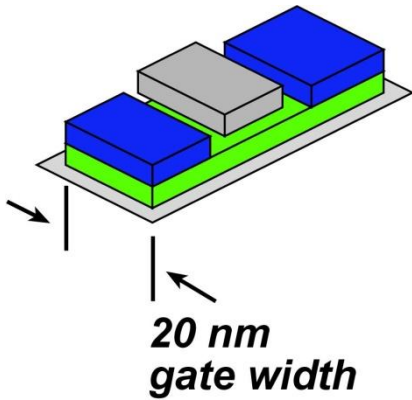
$$C_{total}=2.9\ \text{fF}\ (\text{various multipliers})$$

$$\text{delay}=39\ \text{ps}$$

$$C_{total}V_{DD}^2=0.11\ \text{fJ}$$

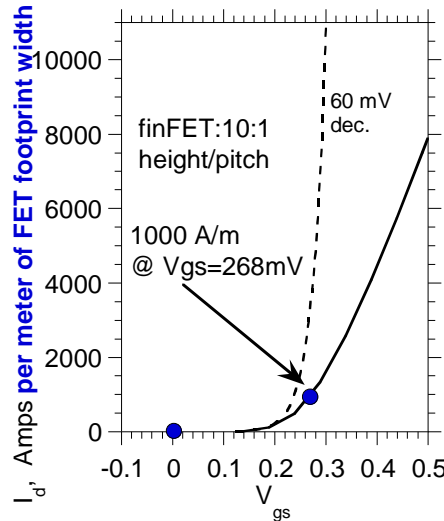
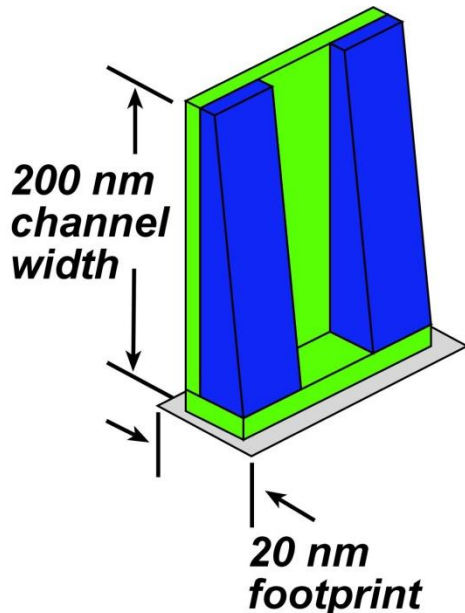
Why tall finFETs ? Why Not Just Subthreshold Logic ?

Planar FET, $V_{dd}=268$ mV



$$I_{on} = 2.0 \mu\text{A}$$

tall finFET, $V_{dd}=268$ mV

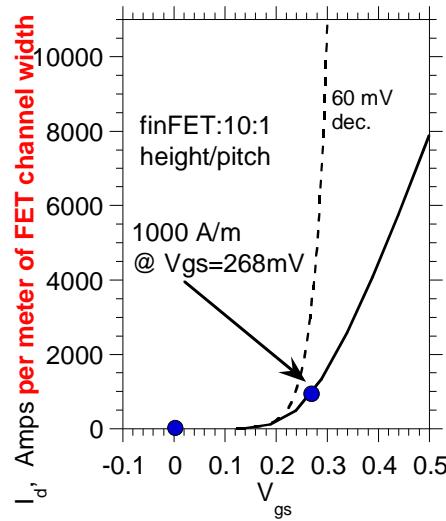
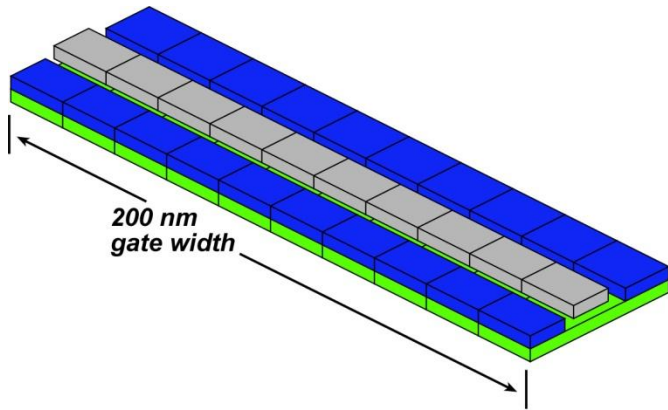


$$I_{on} = 20 \mu\text{A}$$

Low $I_{on} \rightarrow$ large CV_{DD}/I_{on} delay, subthreshold logic is slow.

Why tall finFETs ? Why Not Just Subthreshold Logic ?

Planar FET, $V_{dd}=268\text{ mV}$

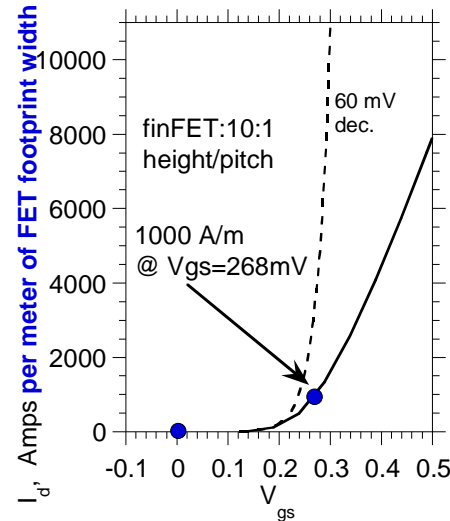
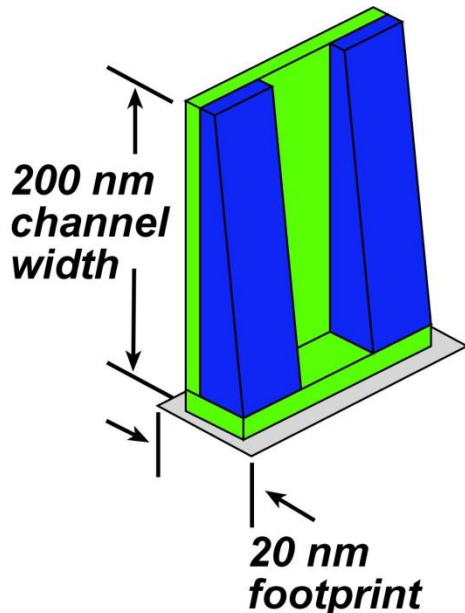


$I_{on}=20\text{ }\mu\text{A}$

**Die size
increased 10:1**

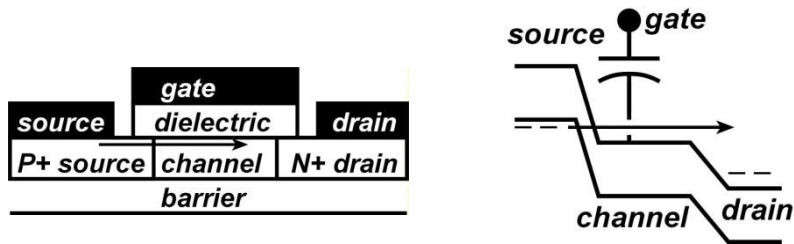
(also: longer interconnects, etc)

tall finFET, $V_{dd}=268\text{ mV}$



$I_{on}=20\text{ }\mu\text{A}$

Tunnel FETs & High-Aspect-Ratio Fins



Quick performance estimate:

Assume, for a moment, that P/N tunneling probability is 10%*.

Typical of the best reported ohmic contacts.*

Then on-currents for tunnel FETs are ~10:1 smaller than that of normal FETs.

Unless I_{on}/W_g is high, tunnel FETs will suffer from either large $C_{wire} V/I$ gate delays or (increasing FET widths) large die areas.

Using high-aspect ratio fin structures, tunnel FET drive currents can be increased. Parasitic fringing capacitance will then also contribute to CV/I & CV^2 .

Contact to N-InGaAs @ $6E19/cm^3$ doping: $m^=0.1m_0$, 0.2 eV, 0.5 nm barrier
Baraskar, et al: Journal of Applied Physics, 114, 154516 (2013)

finFETs Defined by Atomic Layer Epitaxy

Fin thickness defined by Atomic layer epitaxy

→ nm thickness control

Fin height defined by sidewall growth

→ 200 nm high fins

Benefits:

Enables ~4 nm fin bodies → 8 nm gate length

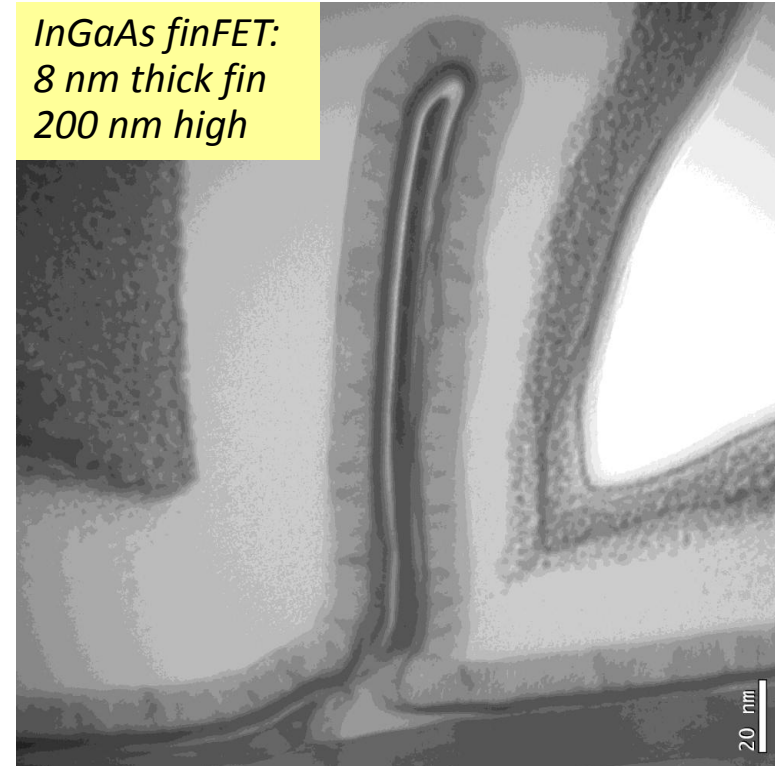
10:1 more current per unit die area

→ **smaller IC die area**

Enables high speed, ultra low-power logic,

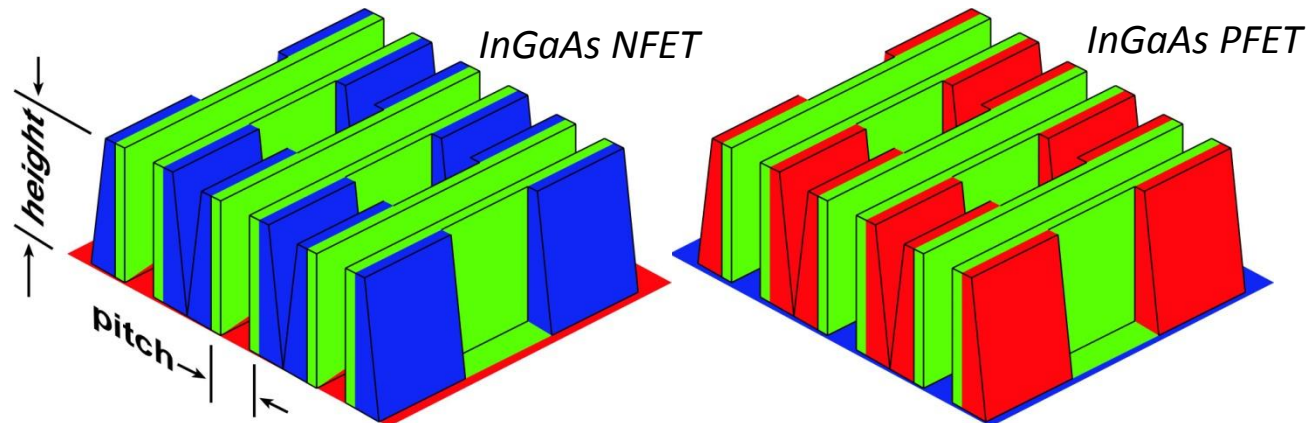
$V_{dd} \sim 300 \text{ mV}$

**InGaAs finFET:
8 nm thick fin
200 nm high**



D. Elias, DRC 2013, June, Notre Dame

height \gg pitch



(end)

Backups

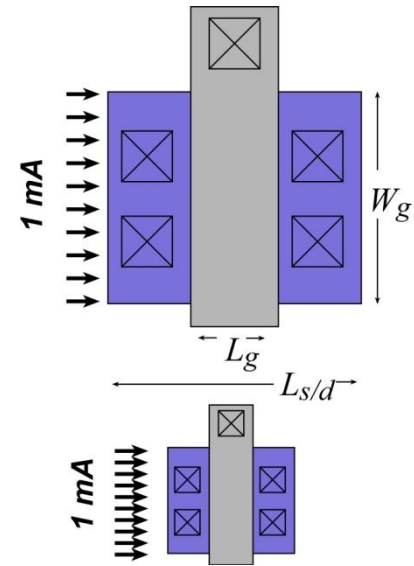
Lithographic Scaling vs. 3-D for High-Density Logic

Past VLSI Scaling: more FETs per IC because of

- (1) shorter gate & contact lengths
- (2) increased mA/micron \rightarrow less gate width W_g .

Today, I_{on}/W_g (mA/micron) is not increasing
 soon, once S/D tunneling dominates, I_{on}/W_g will start to **decrease**

Sub-16nm lithography is also difficult.
 Further reductions in feature size are expensive.

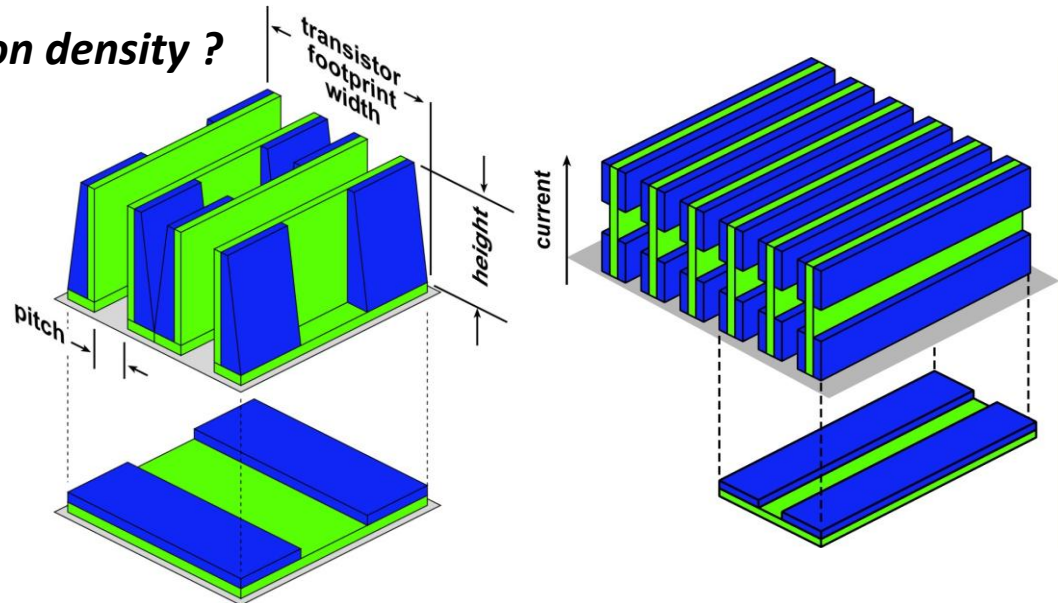


Can 3-D transistors increase integration density ?

Clear: increased I_{on}
 for a given FET footprint size.

Clear: can suppress S/D tunneling:
 tunneling distance \gg
 lithographic distance.

Less clear: decreased size
 of minimum-geometry FET.

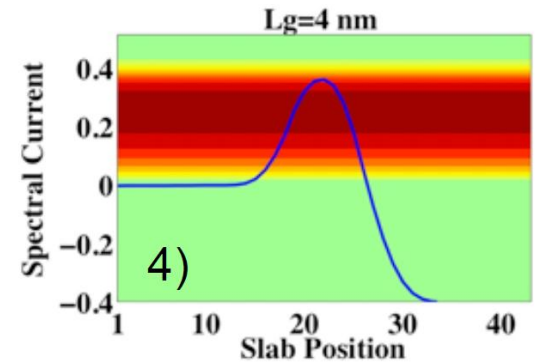


Scaling in the S/D tunneling limit

**At 4-8 nm Gate Lengths,
high leakage from source/drain tunneling**

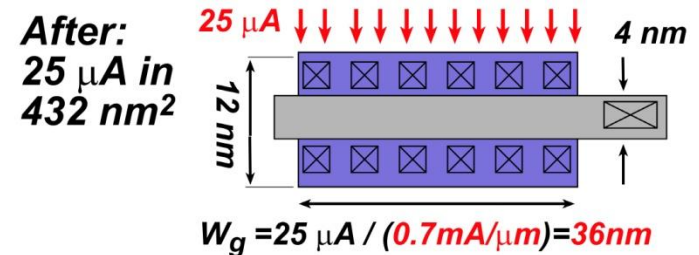
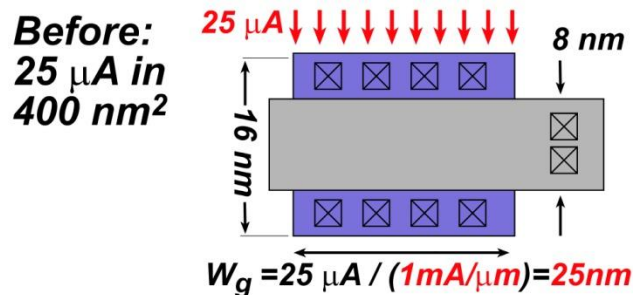
$$J_{S/D\text{tunnel}} \propto \exp\left(-\left(2m^*qV_{th}\right)^{1/2} L_g / \hbar\right)$$

increases exponentially as gate length is reduced.



Reducing tunneling through increased mass can be counterproductive

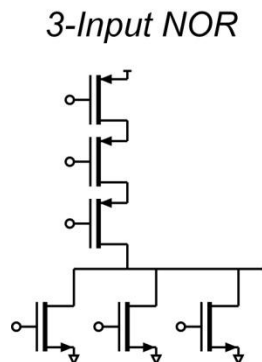
increasing m^ → less tunneling, but lower FET on-current → need more die area*



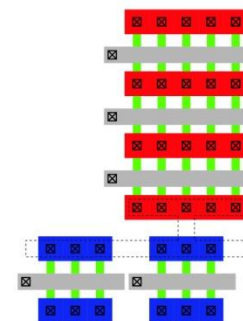
**Instead: Ultra-tall fins to increase
the integration density**

example: 3-input NOR gate.

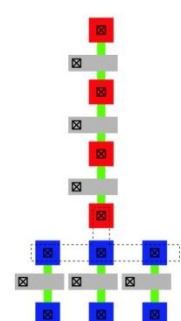
other cases: clock & interconnect drivers



low-current fins

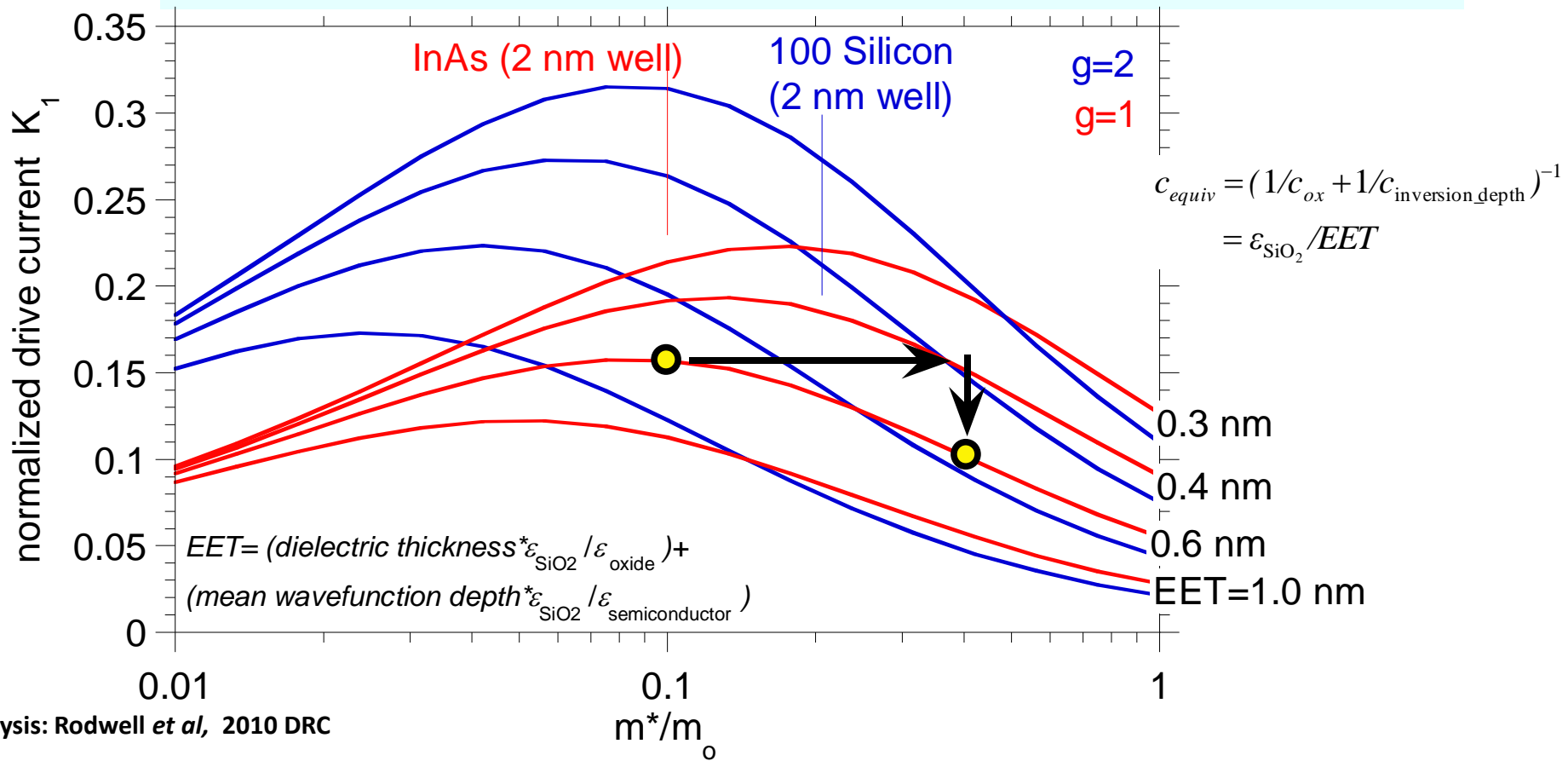


high-current fins



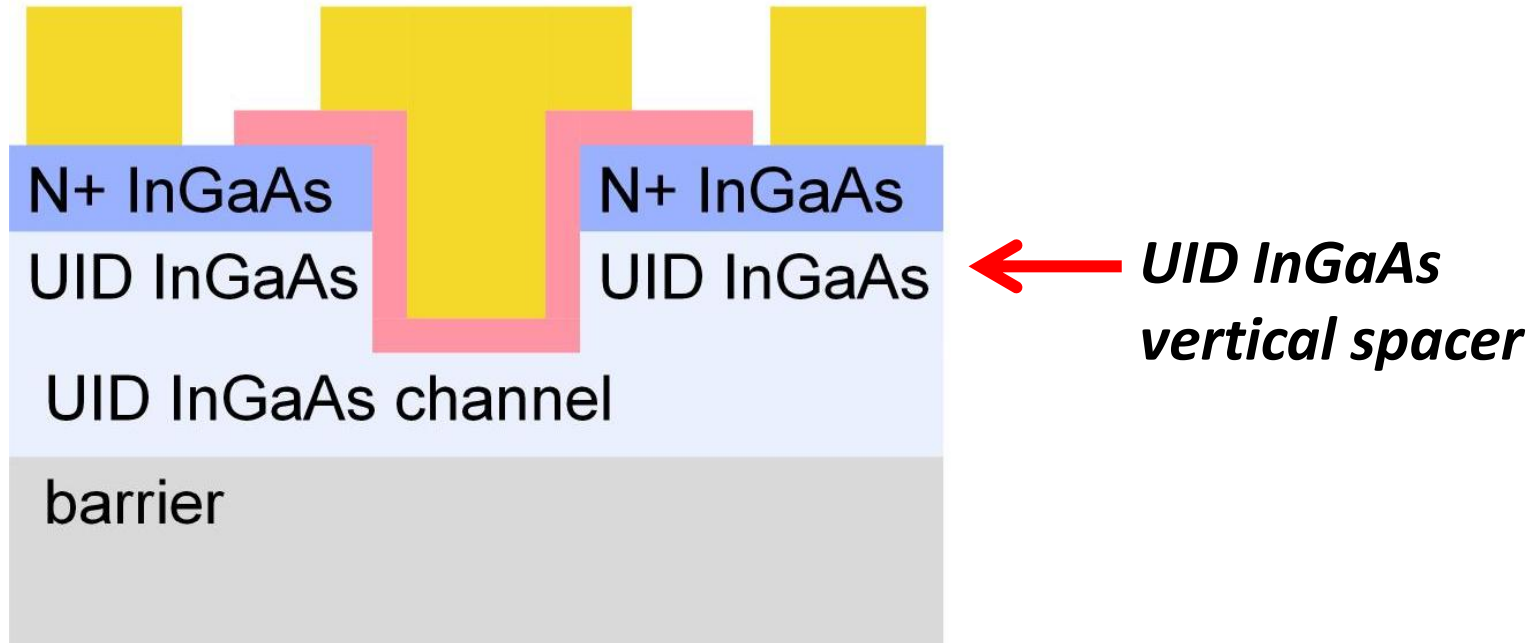
Minimum gate length: source-drain tunneling (2)

$$J = \underline{K_1} \cdot \left(84 \frac{\text{mA}}{\mu\text{m}} \right) \cdot \left(\frac{V_{gs} - V_{th}}{1 \text{ V}} \right)^{3/2}, \quad \text{where } \underline{K_1} = \frac{g \cdot (m^*/m_o)^{1/2}}{\left(1 + (c_{dos,o} / c_{equiv}) \cdot g \cdot (m^*/m_o) \right)^{3/2}}$$



increased $m^* \rightarrow$ decreased $I_{on}/W_g \rightarrow$ decreases packing density

Geometric Solutions to S/D Tunneling



Transport (& tunneling) distance larger than lithographic gate length.

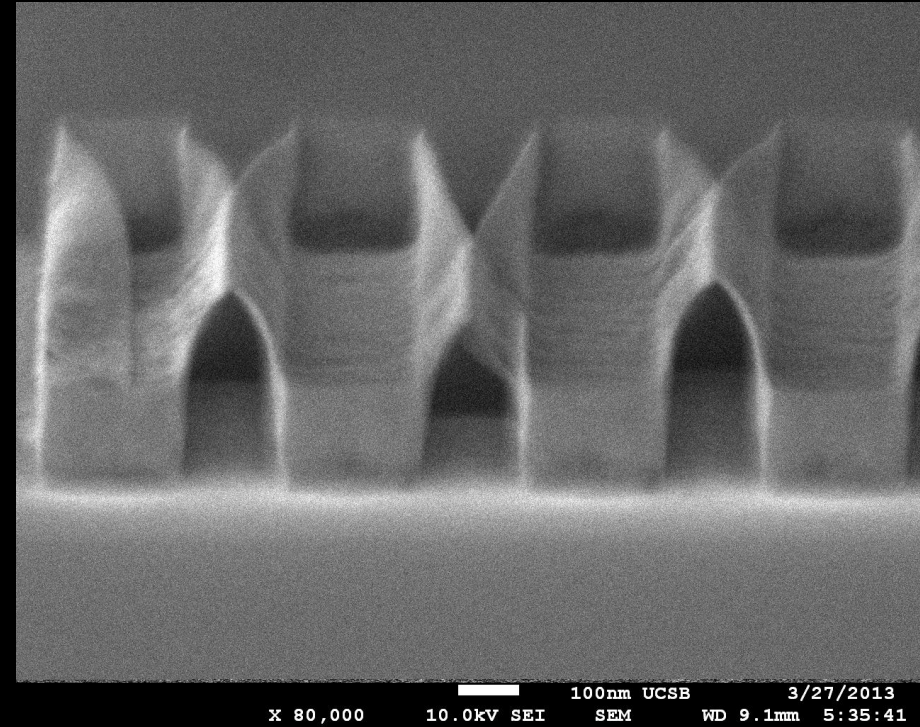
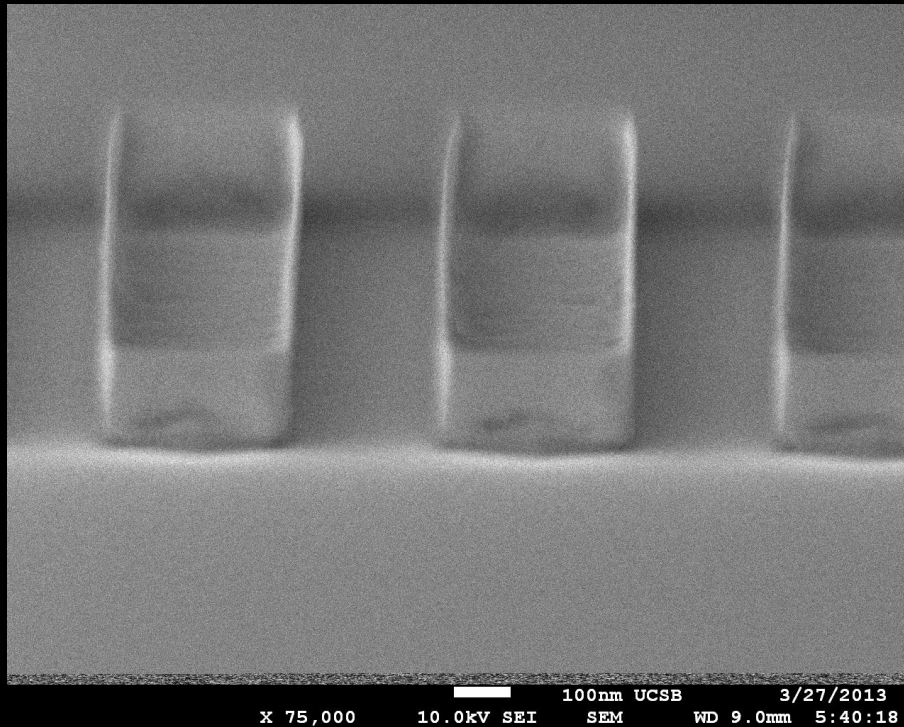
Feature used NOW in our current planar FETs.

Can be incorporated in high-aspect-ratio finFETs

Why Not Release Fins Before S/D Regrowth ?

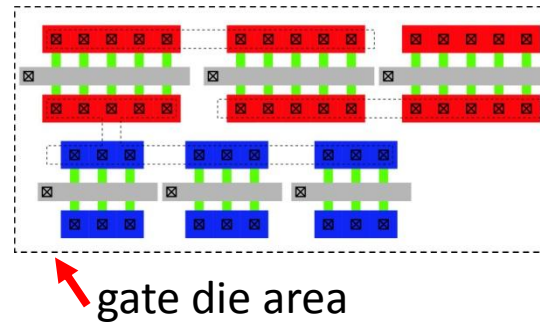
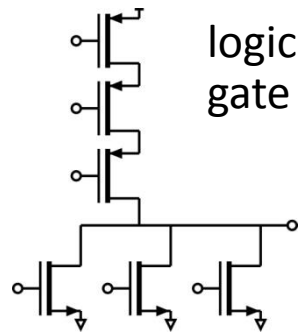
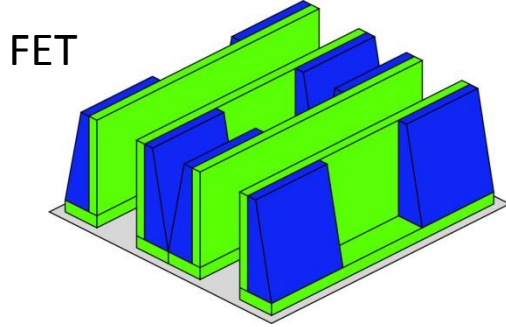
D. Elias, DRC 2013, June, Notre Dame

Images of released ~10 nm fins:

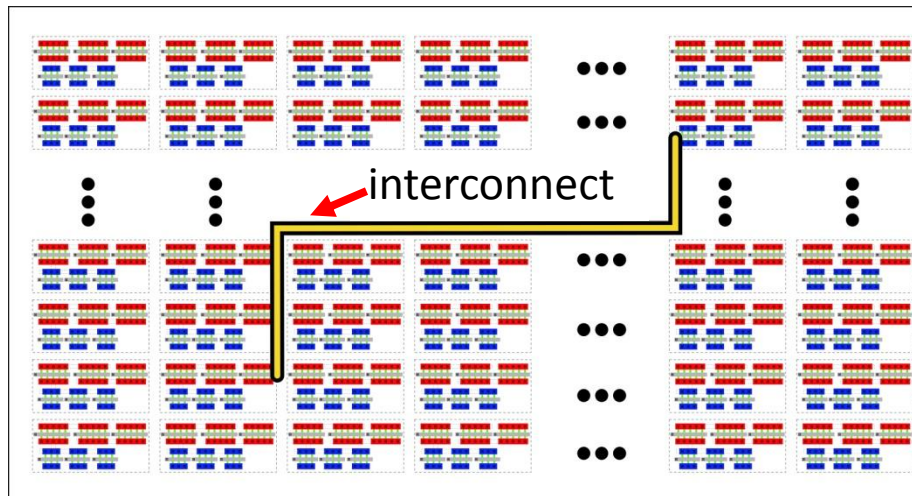


S/D regrowth provides mechanical support

Wire Lengths & Wire Capacitances in VLSI



Integrated Circuit Layout



$$\text{IC area} \propto \text{logic gate area}$$

$$\propto \text{transistor footprint area}$$

$$\text{mean wire length} \propto (\text{transistor footprint area})^{1/2}$$

$$\text{wiring capacitance} \approx 0.2 \text{ fF}/\mu\text{m}$$

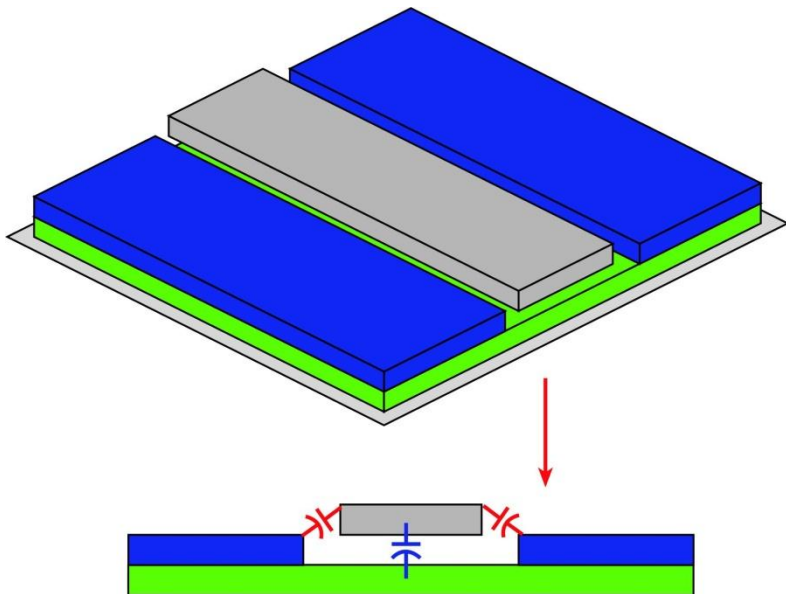
$$\text{mean wiring capacitance}$$

$$\propto (\text{transistor footprint area})^{1/2}$$

more current per fin \rightarrow less fins needed \rightarrow higher integration density

more current per fin \rightarrow shorter wires \rightarrow smaller $C_{\text{wire}} V_{\text{dd}}/I$, $C_{\text{wire}} V_{\text{dd}}^2/2$

FET Capacitances, Interconnect Capacitances



⊥ fringing capacitance

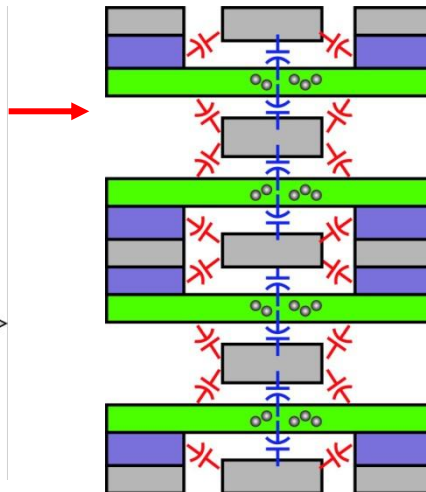
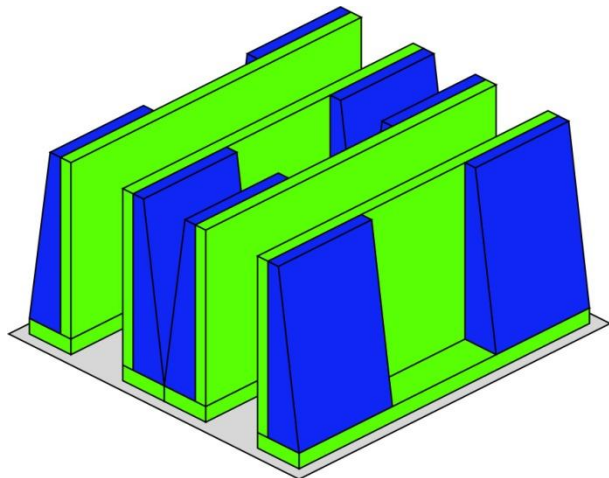
$$C_{gd,f} = C_{gs,f} \approx 0.3 \text{ fF}/\mu\text{m} \cdot W_g$$

⊥ gate - channel capacitance

$$C_{g-ch} \approx I_{\text{on}} L_g / v_{\text{injection}} V_{DD}$$

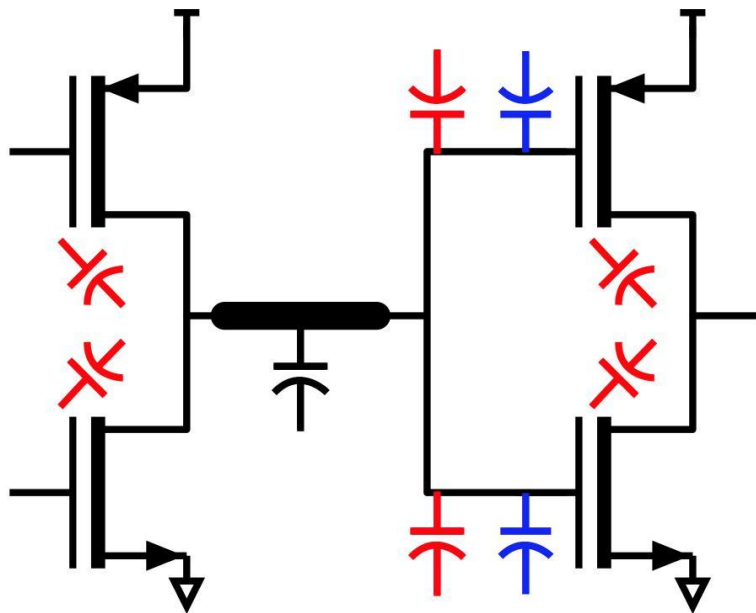
⊥ interconnect capacitance


$$C_{\text{wire}} \approx 0.2 \text{ fF}/\mu\text{m} \cdot L_{\text{wire}}$$





**Similar capacitances
in finFET**

Gate Capacitance, Energy, and Delay



 fringing capacitance
 $C_{gd,f} = C_{gs,f} \approx 0.3 \text{ fF}/\mu\text{m} \cdot W_g$

 gate - channel capacitance
 $C_{g-ch} \approx I_{on} L_g / v_{injection} V_{DD}$

 interconnect capacitance
 $C_{wire} \approx 0.2 \text{ fF}/\mu\text{m} \cdot L_{wire}$

$C_{total} \approx$	$0.2 \text{ fF}/\mu\text{m} \cdot L_{wire}$	interconnect
	$+ 15 \cdot 0.3 \text{ fF}/\mu\text{m} \cdot W_g$	fringing
	$+ 3 \cdot I_{on} L_g / v_{injection} V_{DD}$	gate - channel

Assumes:

- (1) Charge-control analysis*
- (2) $I_{on,PFET} / W_g = 0.5 \cdot I_{on,NFET} / W_g$
- (3) FO=FI=1

$\tau \approx C_{total} V_{DD} / 2I_{on}$ delay

$E_{sw} \approx C_{total} V_{DD}^2 / 2$ energy