Advanced Process Technologies

Prof. Adam Teman

16 March 2022

Special credit to Alvin Loke and Or Nahum for the material and wonderful explanations!

Emerging Nanoscaled Integrated Circuits and Systems Labs





Lecture Overview





Moving to the third dimension





The breakdown of Dennard's Law

• Robert Dennard (1974) observed that continuous scaling is enabled by:

- Scaling the transistor dimensions by 30%
- Scaling the supply voltage by 30%
- This results in:
 - 2X area reduction
 - 30%-40% increase in frequency
 - Constant power density
- However:
 - Voltage wasn't scaled according to the model
 - Frequency was raised faster than the model
 - Leakage, due to VT scaling, wasn't taken into account
- This led to the infamous power crisis



Source: S. Borkar (Intel)

The problem with voltage scaling

- In order to keep constant power density, the supply voltage needs to be scaled
- However, this also requires
 scaling of the threshold voltage
- But the off current of a transistor is exponentially dependent on VT

$$I_{\rm off} [nA] = 100 \frac{W}{L} e^{-V_{\rm T}} / n\phi_{\rm T} = 100 \frac{W}{L} 10^{-V_{\rm T}} / n\phi_{\rm T} =$$

 This is limited by both subthreshold swing and DIBL.



The Multi-Gate Solution

6



Introducing the FinFET



3D Structure

- More I_{on} & gm per area
- Fully depleted channel
 - Less DIBL
 - Less RDF mismatch
 - Negligible body effect
- Quantized channel width
- Problematic Parasitics
 - High S/D resistance
 - High S/D coupling to gate

© Adam Teman, 2022

Planar transistor

FinFET transistor

Main Parameters



- Lower power
- Lower leakage
- Higher intrinsic gain
- Better switch
- Better mismatch
- Smaller area

 Summary: Improved PPA



Fabricating a FinFET





Planar vs. FinFET Fabrication

- We'll start by remembering how planar CMOS is constructed.
 - We patterned the poly gate, oxide and diffusions using a self aligned process
 - Now the backend layers can be made by litho, etch, clean, deposit, polish
- But how do we make Fins above the substrate plane?
 - We carve them out of the substrate!
 - And then cover them with oxide and poly.







© Adam Teman, 2022



Fin and Gate Patterning

 The first critical step of the FinFET process is to make very thin fins at a very tight pitch and build tightly pitched gates.

- Immersion Litho
- Phase shifting
- Double patterning (a.k.a. LE-LE)
- Self-aligned
 Patterning
 (SADP/SAQP)

Double Patterning (LE-LE)

Introduce "colored" masks for pitch splitting

Limited by misalignment between masks.
 Minimum Pitch

Fins and Gates Source/Drain IVILUL DLUL

Self-Aligned Double Patterning (SADP)

- We can achieve the same resolution using a *self-aligned* technique
 - Start by constructing "mandrels" with the minimum litho pitch.
 - Next, create spacers around the mandrels.
 - Remove the hard mask and the spacers are at half the pitch.
 - Note that Line Edge Roughness is correlated, reducing L variation!

Self-Aligned Quadruple Patterning (SADP)

- And just repeat the process to double the resolution again.
 - Just remember that this cost us another few masks and many steps.

Additional Advanced Litho Points

- You can only print one line width!
 - To get multiple line widths, change mandrel width and spacing.
- Each multi-patterned layer is unidirectional.
 - No more wrong-way routes or jogs!
- Use orthogonal cut masks to break patterns.
 - Orthogonal cut mask according to CD.
 - Eliminate corner rounding on fins/gates.

block mask

Woo, et al. Globalfoundries

Auth, et al. Intel

cut mask

pattern

Source/Drain Regions

- We have seen how to make the "array" of fins and gates.
- But we have a few problems with the S/D regions:
 - 1. They are really small and delicate \rightarrow hard to contact to.
 - 2. They have high resistance.
 - 3. We need to apply stressors to improve mobility.
- Therefore, we build an epitaxial area on the fins

A note about Stress

- A silicon channel is piezoresistive
 - Lattice strain affects mobility.
 - Tensile stress improves NMOS, compressive stress for PMOS
 - Depends on lattice orientation: (100) fin top vs. (110) fin sidewall

Stress has been more effective for PMOS[®]

- This has caused beta (N/P) ratio to fall to about unity at 7nm.
- Expected to change for nanosheets

Hi-K <u>- MG</u>

High-K – Metal Gate (HKMG)

• Two problems with the transistor gate:

- 1. Thinning Oxide → Gate Leakage
- Polysilicon → high resistance, poly diffusion
- Therefore, at 45nm-28nm move to HKMG
 - High-K dielectrics enable thicker oxides
 - Metal gates improve resistance, EOT
- However:
 - High-K Materials have lower Energy Bandgap
 - Metal gates are sensitive to high temperatures during annealing

Gate dielectric Material	Dielectric constant (k)	Energy bandgap Eg (eV)	Conduction band offset ΔEc (eV)	Valence band offset ΔEc(eV)
SiO ₂	3.9	9	3.5	4.4
Al_2O_3	8	8.8	3	4.7
TiO ₂	80	3.5	1.1	1.3
ZrO_2	25	5.8	1.4	3.3
HfO_2	25	5.8	1.4	3.3

Touati, et al., J. New Technol. Mater.

19

Hi-K - MG

MEOL

Metal Gate and VT adjustment

• VT of a transistor is primarily set by:

- The workfunction difference between gate and substrate
- The doping in the junction
- The backgate capacitance
- In older technologies
 - VT adjustment was achieved through channel implants
 - However, random doping fluctuations caused huge variation

Hi-K - MG

- FinFETs have intrinsic channels
 - Therefore, RDF has basically "gone away"
 - VT adjustment is done through the workfunction of the metal gate

Replacement Metal Gate (RMG)

- Metal Gates are sensitive to high-temperature process steps required for S/D engineering (epi, stress, etc.)
- Therefore, a "gate-last" approach is used:
 - Form the gates with polysilicon. Also known as a dummy gate.
 - After S/D formation, etch poly gates.
 - Partially fill in with barrier and workfunction metal.
 - Fill in with low resistance metal.
- This also impairs the Silicide, increasing S/D resistance

21

Hi-K <u>- MG</u>

MEOL

BEOL

© Adam Teman, 2022

Middle-end-of-the-line (MEOL)

- In older technologies, making contacts was easy
 - The contacted gate pitch (CGP) was large enough to place gates and contacts next to each other.
 - But at tight gate pitches, any misalignment can short the contact to the gate.
- Therefore, Self-Aligned Contacts (SAC) were introduced
 - A dielectric cap is added on top of the gate so that if the contact overlaps the gate, no short occurs.

22

Contac

Contac

Middle-end-of-the-line (MEOL)

- SAC process is more complex in FinFET
 - RMG is harder to recess than gate-last approach
- Contact Over Active Gate (COAG) is desired
 - Introduced by Intel in 10nm process
- The MEOL process now includes many steps
 - Cap gate and create SAC
 - Cap contact and create COAG
 - Create additional VIA0 to go through caps
- **Disadvantages**:
 - High Gate to S/D Contact capacitance
 - High S/D, MEOL & lower BEOL resistance

23

Hi-K - MG

MEOL

Backend-of-the-line (BEOL)

- Copper interconnect replaced aluminum
 - Lower resistance, Improved electromigration
 - Dual Damascene Process
- However, a barrier and liner are required
 - Plated copper area is reduced, increasing resistance
- At 10nm Intel started using Cobalt on M0 & M1
 - 5-10X Electromigration, 2X resistance

Jacob, et al., Globalfoundries

Intel 10nm metal interconnects cross section wikichip

© Adam Teman. 2022

BEOL

FinFET Layout

Layout Complexity

© Adam Teman, 2022

Planar vs. FinFET

FinFET Stack (5nm)

Diffusion Breaks

• A diffusion break is required between two active areas:

- Blocks Epitaxial growth
- Provides "back wall" for stressors
- "Double Diffusion Break" (DDB) is done with two dummy gates and STI
- This wastes a lot of area
 - Instead use only one dummy gate
 - "Single diffusion break" (SDB) saves area but is complex to process

Density & Floorplan Considerations

- Critical process steps are extremely sensitive to pattern density & loading
 - 1000s of DRCs, many very tough to pass, increasingly restrictive & foreign
- DRCs reduce unmodeled long-range systematic & random variation
 - → iterative rework of smaller cells
 - Area, perimeter, gradient
 - Contacts, vias, cuts, tight-pitch metal
 - Larger checking windows
 - Density union of multiple metal levels
- Floorplanning more tedious & bloated
 - More dummy gates, well taps, guard rings
 - Wasteful transitions between different device types & pattern densities

Yang, et al., Qualcomm © Adam Teman, 2022

Layout Dependent Effects and Parasitics

Layout Dependent Effects (LDEs)

- Fabricated device characteristics have shown a high dependency on layout features for several generations:
 - Well Proximity Effects (WPE)
 - Length of Diffusion (LOD)
 - Oxide-to-Oxide Spacing (OSE)
- FinFETs further introduce LDEs:
 - Stress LDEs more significant due to stronger stressors
 - Gate cut stress
 - HKMG LDEs
- All of this causes more pre- to post-layout simulation differences

Stress LDEs in FinFETs

- Stress LDEs are caused by:
 - Longer diffusions (OD Length)
 - Wider diffusions (I_D /fin not constant vs. # fins)
 - Oxide spacing
 - Gate pitch
- Models capture $\Delta \mu \& \Delta V_T$ (some effects as early as 130nm)

Gate Cut Stress LDE

- Gate cut disrupts mechanical support of continuous gate & stress near cut
 - $\rightarrow \Delta \mu \& \Delta V_T$, modeled in post-layout netlist starting in 16/14nm

HKMG LDEs

Metal Boundary Effect

- $\Delta V_{\rm T}$ near border of different $\Phi_{\rm M}$ due to Interdiffusion of $\Phi_{\rm M}$.
- Mitigated with gate cut but costs area
- Models capture $\Delta \mu \& \Delta V_{\rm T}$.

Density Gradient Effect (DGE)

- Gate density gradients
 - $\rightarrow \Delta V_{\rm T}$ & variation from RMG CMP dishing
- $\Phi_{
 m M}$ influenced by metal fill & sidewall $\Phi_{
 m M}$
- Not modeled, contained with DRC

Process Loading Variation

- Local pattern density modulates deposition rate, etch rate/profile & CD
 - Deposition loading: Spacer width variation (gate and metal CD)
 - Epitaxy loading: S/D volume variation (S/D resistance & channel stress)
 - Etch loading: Depth/profile variation (Lgate, fin & metal height)
 - Rapid Thermal Annealing (device variation)
 - Chemical Mechanical Polishing (variation in STI, poly, RMG, MEOL & BEOL)

Parasitic Resistance and Capacitance

High Resistance:

- Contacts, Metal Gate, Low Metals
- High Capacitance:
 - Tight metal pitches
 - S/D trench contacts & gate form vertical plate capacitors

Сар	Planar	FinFET		
C _f		2.0X		
C _{co}	1X	1.8X		
C _{gb}		1.2X		
[Lawrence Loh. ISSCC'18]				

Res	Planar	FinFET
R _g		1.5-2.5X
R _{ct}	1.V	1.5X
R _{int}	IX	1.3X
R _{diff}		1.3X

© Adam Teman, 2022

Gate Resistance

Metal gate should have lower resistance than poly gate

- But gates are very thin.
- Workfunction Metals have high resistance.
- Lower resistance metal filled on top of workfunction metal
 - But still, this is very little conductive metal

Another point

 It's tough to make thick oxide I/Os!

Wu & Chan, HKUST / Lee, Intel

Diffusion & MEOL Resistance

- Traditionally:
 - Try to share diffusions for smaller S/D capacitance
 - But S/D resistance is now a much bigger problem than capacitance
- This becomes challenging for high-current circuits
 - e.g., I/O drivers, clock buffers

BEOL Resistance

Copper interconnect used for low resistivity

- However, copper diffuses into ILD.
- Need barrier, liner and seed layers.
- Local interconnect (MX) aggressive pitch scaling
 - Dense logic routing→less die area & cost
 - But not much copper left in the wire...
 - 6X rise in resistance from 80nm to 48nm pitch!
- Therefore:
 - Remove seed, liner layers.
 - Use cobalt & ruthenium which don't require a barrier
 - Despite higher ρ less wire resistance!
 - Use M3 and up for inter-cell routing.

Cu

Parasitics Summary...

Half of your performance from scaling is going here!

Courtesy of Greg Yeric, ARMTechCon 2016

https://madematics.com/2011/09/13/nanotoilet/

© Adam Teman, 2022

FinFET Node Models

• FET models

- BSIM-CMG based on channel surface potential, less equation fitting
- *Target-based* for latest nodes, more silicon influence in mature nodes
- Prone to model-vs.-silicon gap from increasing density & loading effects

BEOL models

- Electrical information provided, limited to no physical stack-up details
- Less pessimistic corners for relaxed timing closure (customer pressure)
- Usual reliability models (HCI, BTI, TDDB, EM)
 - Vague allowable VDD, depends on application
- Foundries extremely paranoid to protect their technology IP from competitors
 - Process corner methodologies & many model parameters encrypted
 - Limited physical information available even basic dimensions (e.g., L_{gate}) not real
 - CD bias & mask booleans to conceal process details (e.g., RMG flow for multiple VT)

Overcoming Process/Model Immaturity

• "All models are wrong, but some are useful"

George Box (1919-2013), British Statistician

• Corollary: "All simulations are wrong, but some are useful"

Alvin Loke, NXP/TSMC

Layout Guideline	Reduced Exposure	
Use continuous OD stress plateau	Stress LDEs (S/D epitaxy, STI)	
Attach dummy FETs to OD ends		
Avoid single-diffusion break		
Use only one ${\it P}_{\rm M}$ in each gate	Metal boundary LDE	
Avoid using gate as interconnect	Gate cut LDE	
Contact gate on both sides	Gate resistance, Φ_{M} tuning to	
Use groups of fewer fins	adjust V _T	
Use double-source layout for high- <i>I</i> nets	S/D contact resistance & epitaxy	
Extend S/D contact to land extra via	MEOL & S/D resistance	
Do not push DRCs to limit	DRC updates	

Some Current Trends

A note about process node naming

Scaling rate slower than 0.7x per node

- Node name just marketing number for equivalent PPA
- Physical gate lengths haven't scaled below ~14nm
- 193i Single exposure pitch limit ~80nm
- "Intra Nodes" on official scaling roadmaps
 - Process optimizations for performance improvements and yield enhancements

Zheng Guo et al. Intel

22nm Process 1st Generation FinFET

2nd Generation FinFET

10nm Process 3rd Generation FinFET

2008 2010 2012 2014 2016 2018 2020 Loke, BCICTS

47

Squeezing Out What's Left in FinFETs

Really tough after 4 FinFET generations

- Realistically, never been any low hanging fruit with each new node
- Process innovations & complexity for only incremental gain
- +5% ring oscillator frequency is a big deal

Areas of development (no stones unturned)

- Short-channel control \rightarrow narrower fins, tradeoff vs. μ reduction
- Channel mobility \rightarrow high- μ fin material, e.g., TSMC 5nm
- EOT \rightarrow higher K, thinner & reliable gate dielectrics
- Device variation → fin uniformity & geometry control
- Volumeless VT tuning using only HK dipoles
- $R_{\text{contact}} \rightarrow \text{contact resistivity (interface quality) & area$
- *R*_{gate} → selective bottom-up HKMG deposition
- $C_{GS} \& C_{GD} \rightarrow$ gate spacer K, air gap spacers
- MEOL & BEOL resistance \rightarrow metal resistivity, R_{via}

Vertical Scaling and Fin Depopulation

Migrating to Gate-All-Around (GAA)

- FinFET has poor short-channel control with further L_{gate} scaling
 - Need better short-channel control & more W_{eff} per die area
- Stacked GAA nanowires & nanosheets are promising
 - Nanowires offer better SCE, nanosheets offer better area scaling

stacked nanosheet demonstration

FinFET → Nanosheet

- Current mostly flows along (100) surface
 - Electrons ☺, Holes ☺→NMOS again stronger
- Need inner spacers to reduce gate-to-S/D capacitance
- Disable parasitic planar FET below nanosheets

Backside power delivery

• Burying the power rails *under* the transistors

- Power rails do not take up routing resources
- →Reduced design area
- →Improved IR drop and voltage droops
- →Reduced power.

Intel's PowerVia

backside

power delivery

© Adam Teman. 2022

Conclusions

- "Moore's Law is well and alive. It's not slowing down. It's not even sick" Phillip Wong, TSCM, HotChips 2019
- SoC area scaling now driven primarily by device innovation & DTCO, less by feature size reduction.
- Understand & exploit technology for maximum PPA benefit & efficient design productivity

Wong, TSMC, HotChips, 2019 © Adam Teman, 2022

Main References

- Alvin Loke, et al., "Nanoscale FinFET Technology for Circuit Designers", 2019 CICC Education Sessions, 2020 BCICTS, 2021 MTT-SCV
- Or Nahum, "FinFET Process Overview", 2021
- www.halbleiter.org "Semiconductor Technology from A to Z"
- Jacob, et al., "Scaling Challenges for Advanced CMOS Devices" Globalfoundries