

NAND Device Physics

Akira Goda
Krishna Parat

Acknowledgement:

- Intel and Micron Global NAND Team

Outline

- Introduction – 2D & 3D Cell and Array Architecture
- Read – Cell I-V, Sensing, Interference, MLC, TLC
- Programming – Program Physics, Algorithms, PVS
SSPC, Program Noise, Inhibit, Inhibit schemes
- Erase – Erase Physics, Algorithms
- Read Window Budget – Key components of RWB
- Cell Reliability Overview – Cycling, ICL, SBCL, Read Disturb, Floating Body effects, FPR, ECC

General Schedule

Day 1:

• Intro -	~30 slides	~1.5hr
• Read -	~30 slides	~1.5hr
• Single Cell Programming –	~15 slides	~1.0hr

Day 2:

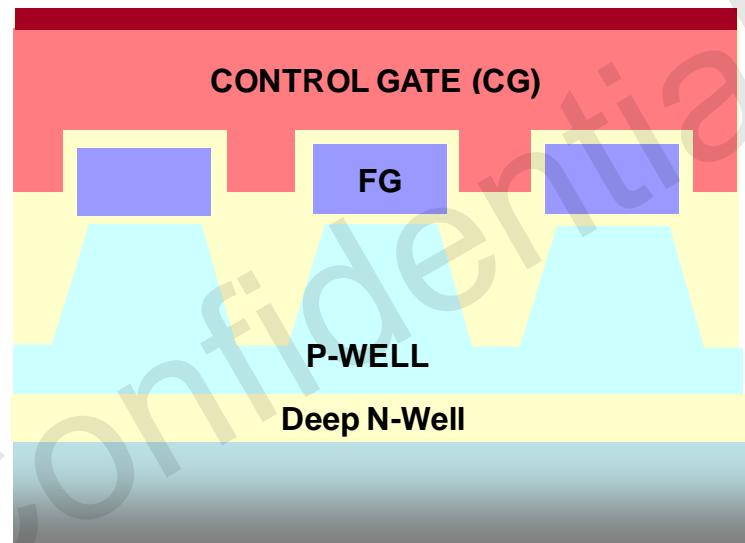
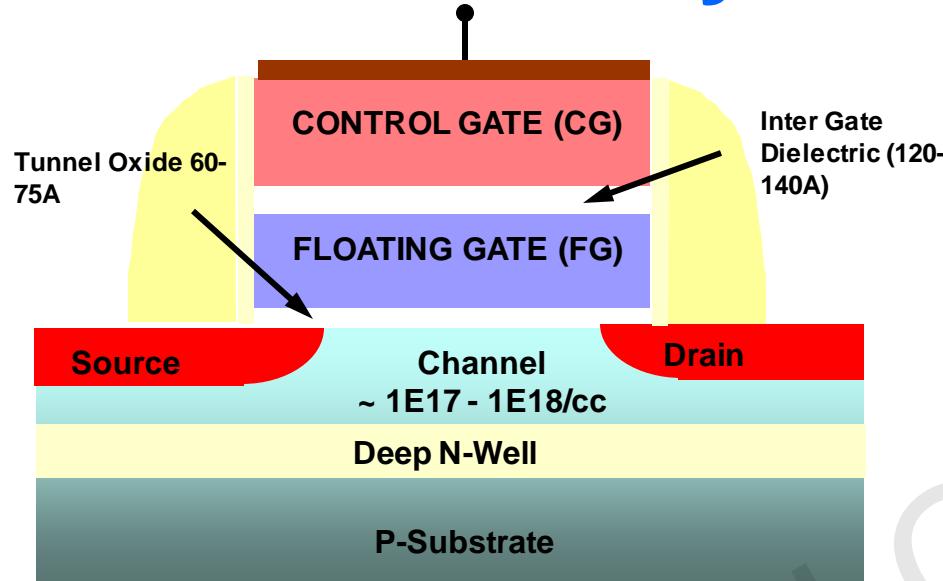
• Array Program & Inhibit –	~30 slides	~1.5hr
• Erase –	~15 slides	~1.0hr
• RWB & Rel –	~25 slides	~1.5hr

➤ Feel Free to Interrupt and ask questions anytime during the class

Introduction

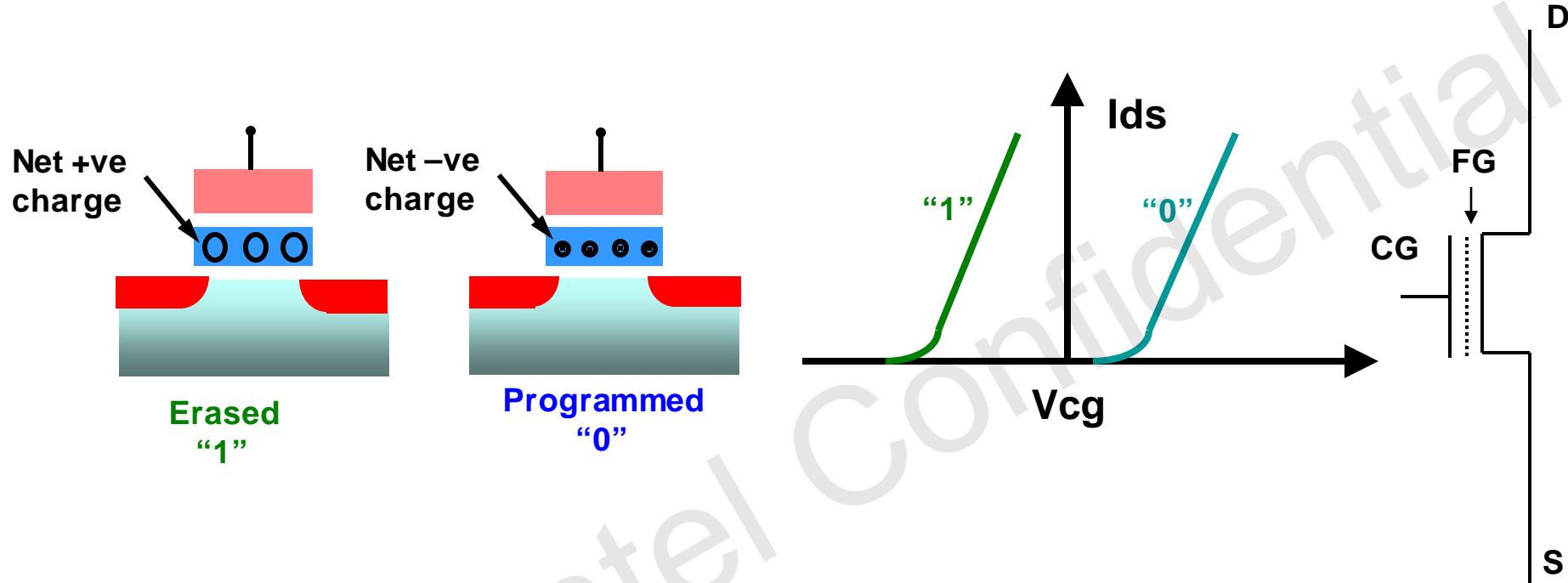
- NAND Flash Memory Basics
- Device structure
- Device operation
- 2D Cell and String
- 2D Cell Structures
- 3D Cell
- 3D Cell Structures
- 2D Array
- 3D Array

Flash Memory Device (2D)



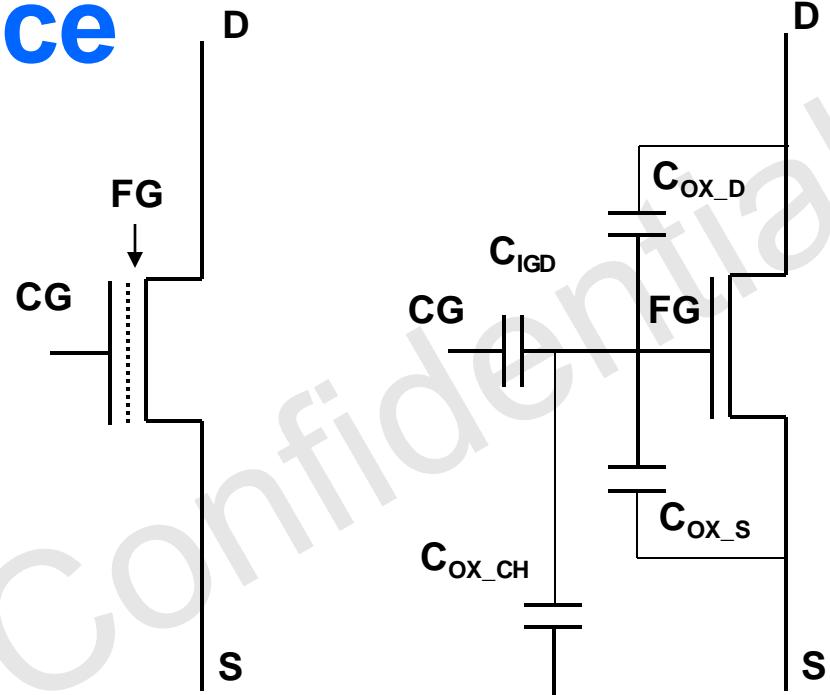
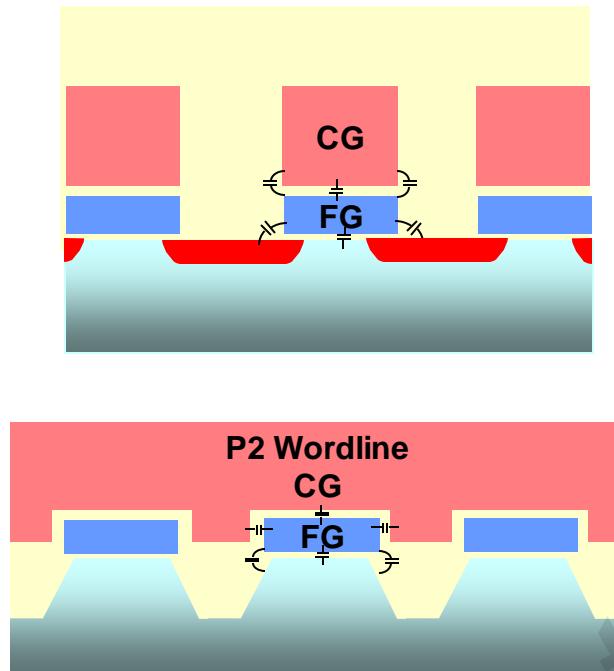
- Stacked Gate NMOS Transistor
 - Floating Gate: Charge storage node
 - Control Gate: Access/Control gate for the transistor
 - Tunnel-oxide: Gate oxide and dielectric through which charge is tunneled in and out of the floating gate (charge storage node)
 - Inter Gate Dielectric: Multi-layer/barrier dielectric with low leakage
 - Lightly doped Source/Drain Junctions optimized for leakage
 - Channel isolated from Si substrate to allow applying voltage to Channel

Flash Memory Device



- Threshold Voltage shift = $-\Delta Q_{FG}/C_{IGD}$ $\left\{ \begin{array}{l} \Delta Q_{FG} = \Delta \text{ in charge on FG} \\ C_{IGD} = \text{CG to FG Capacitance} \end{array} \right.$
- Program = Electrons Stored on the Floating Gate \rightarrow High V_t
- Erase = Remove electrons from the Floating Gate \rightarrow Low V_t
- Read = Look for current through the cell at given gate bias

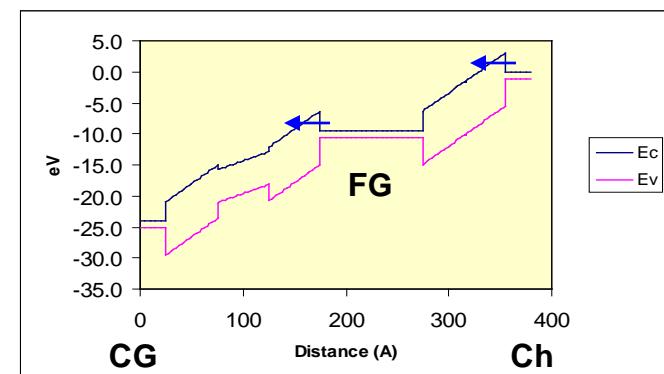
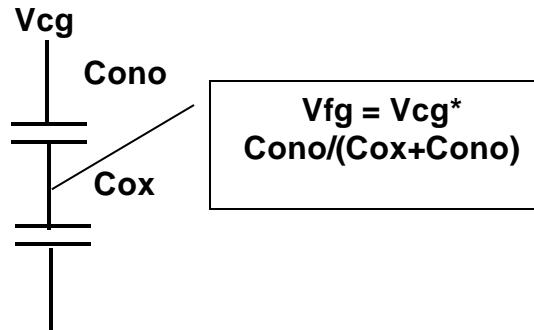
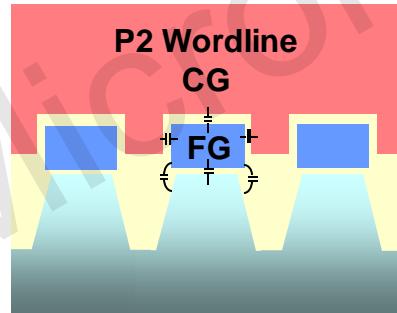
Flash Memory Device



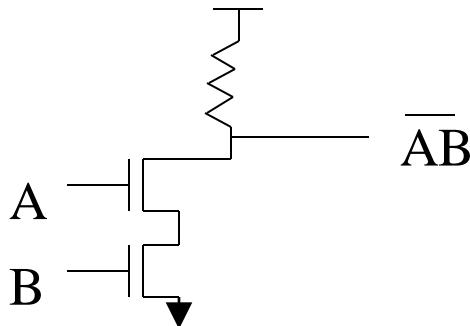
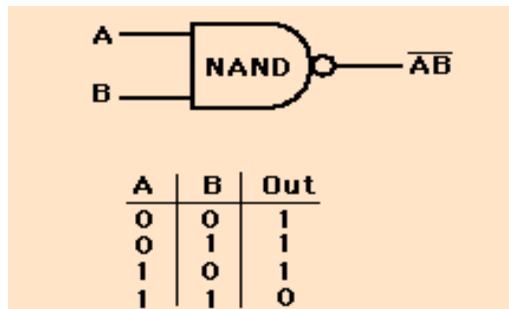
- Cell is a linear capacitor network with n-ch transistor attached
- Various cell capacitances can be calculated using TCAD
- FG Voltage, $V_{FG} = Q_{FG}/C_{TOT} + C_{IGD}/C_{TOT} * V_{CG} + C_{ox_D}/C_{TOT} * V_D + \dots$
 - Where, Q_{FG} is the net charge on FG, and $C_{TOT} = C_{IGD} + C_{ox_D} + \dots$
- Gate Coupling Ratio (GCR) = $dV_{FG}/dV_{CG} = C_{IGD}/C_{TOT}$
- High GCR is desirable for CG to have a good control of FG & Channel

Flash Dielectric Constraints

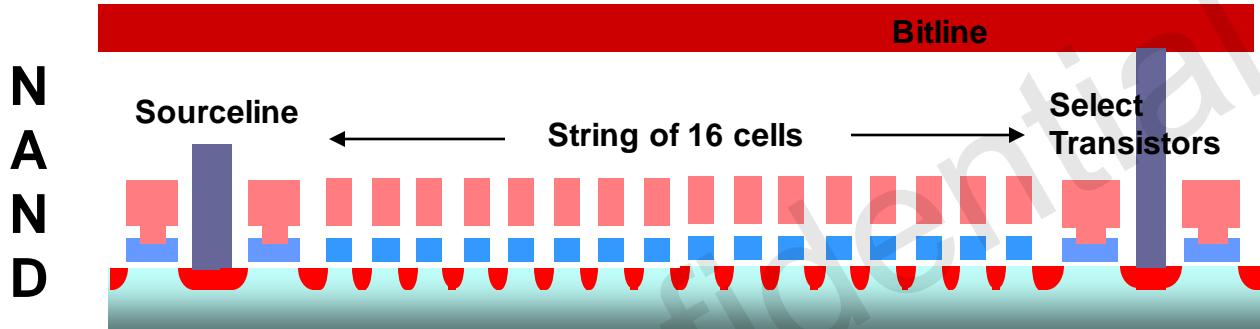
- Data retention limits the Tunnel-oxide scaling to ~55-75A
 - Intrinsically a ~55A silicon dioxide with ~3eV barrier to N-Poly can provide the needed retention (<1E-16A/cm² of leakage).
- IGD thickness dictated by retention as well as the need for large Program/Erase window (high field leakage during Program/Erase)
 - ONO as IGD scaled down to ~110-120A Physical (EOT ~ 100A)
 - Hi-K IGD scales down to ~150-160A Physical (EOT ~ 80A)
 - High barrier important for low leakage
 - Low electrical thickness important for good gate control
 - Good Gate control can be achieved by increasing the coupling area between the CG and FG → CG wrap of the FG
- High doping of the FG and CG desired to minimize poly depletion



NAND Schematic Cross-sections



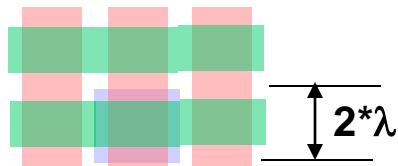
1. Transistors connected thru each other to output & ground.
2. Current flows simultaneously through all transistors



NAND string typically has 32-128 cells (+ Dummy cells). Schematic shows example with 16 cells.

- Individual NAND cells do NOT require a contact → very compact cell layout (2D NAND - typically $4-5\lambda^2$)
- One select transistor (SGD & SGS) at each end of the NAND string
- Select Gates, however, adds additional area overhead of about ~20%

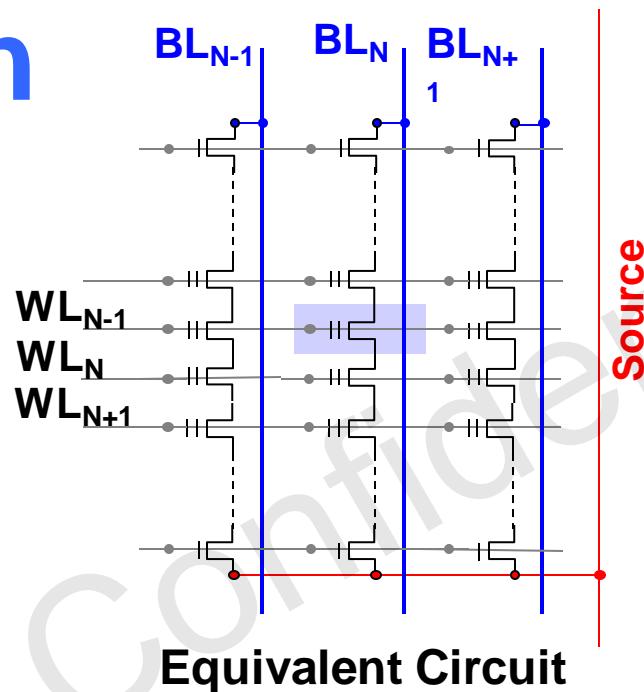
2D NAND Flash



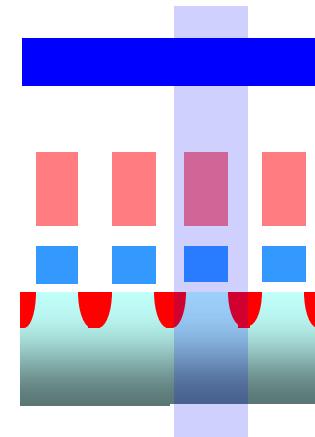
Layout



Y Cross-section



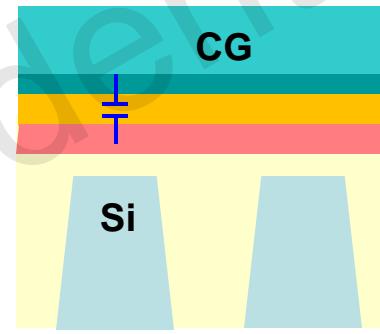
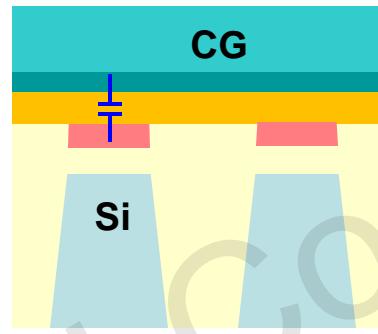
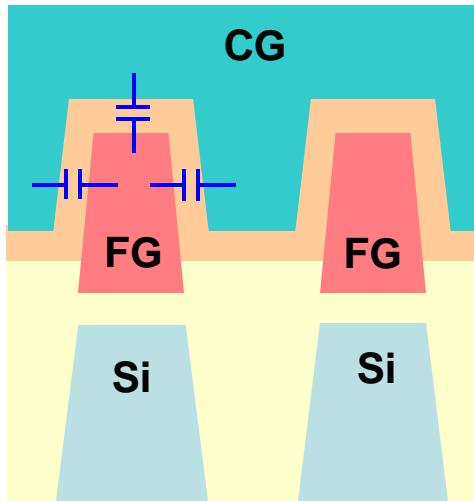
Equivalent Circuit



X Cross-section

Select Transistors shown in the equiv ckt, but not shown in the layout or cross-sections

2D Flash Memory Cell Structures



Wrap FG Cell

Planar Floating Gate Cell

Planar Charge Trap Cell

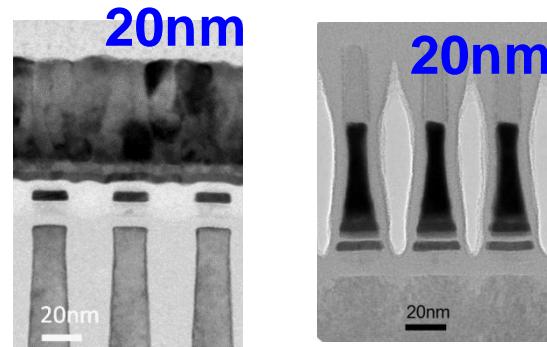
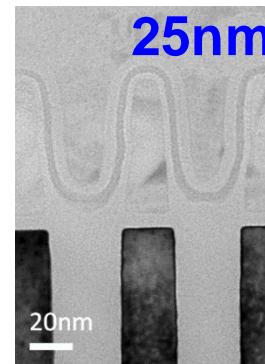
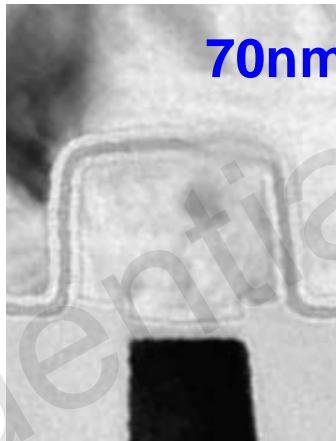
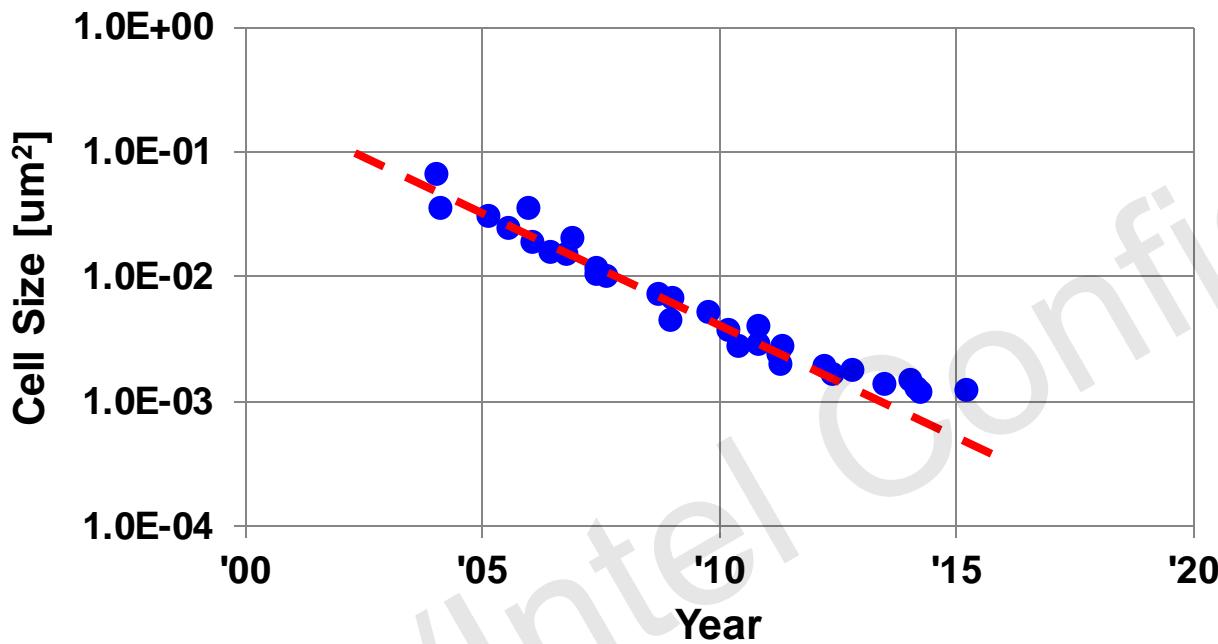
- NAND Technology up to 70s
- Polysilicon Floating Gate
- ONO IPD (IGD)
- Polycide/Metal Control Gate

- 80s & 90s Technology
- Hybrid Floating Gate
- Multi-barrier dielectric Blocking dielectric (BD)
- Metal Control Gate

- SONOS/TANOS
- Nitride / Dielectric Charge storage
- Oxide/High-K Blocking dielectric (BD)
- Metal Control Gate

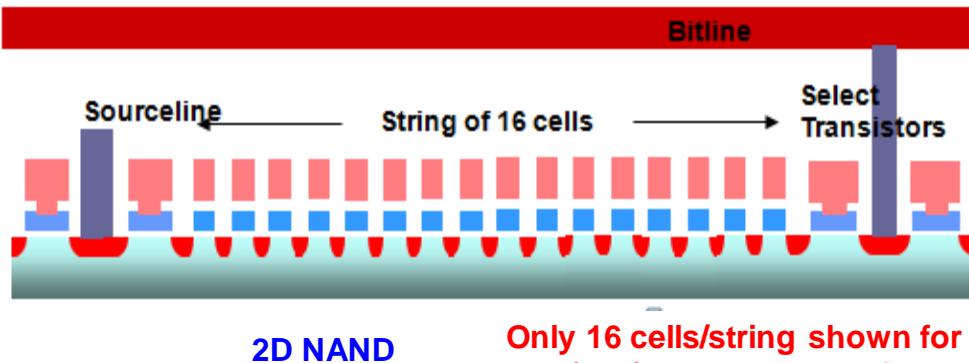
2D NAND Cell Scaling

2D NAND Cell Size Scaling

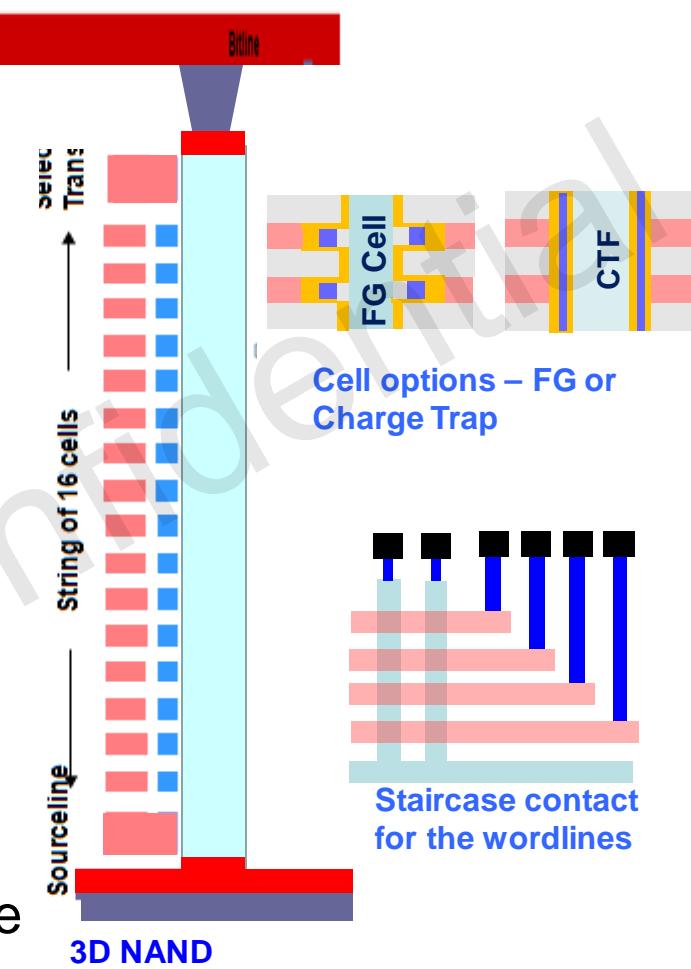


- 2D NAND scaling has slowed down:
 - Lithography Limitations
 - Small Cell Area Effects: Number Fluctuation
 - Proximity Effects: Cell to Cell interference
 - High Electric Field Effects

2D NAND → 3D NAND



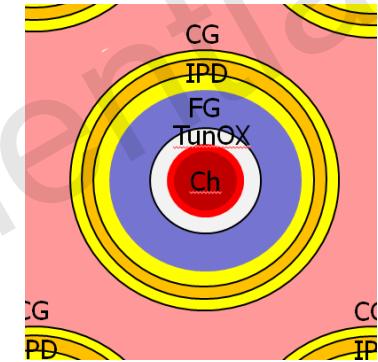
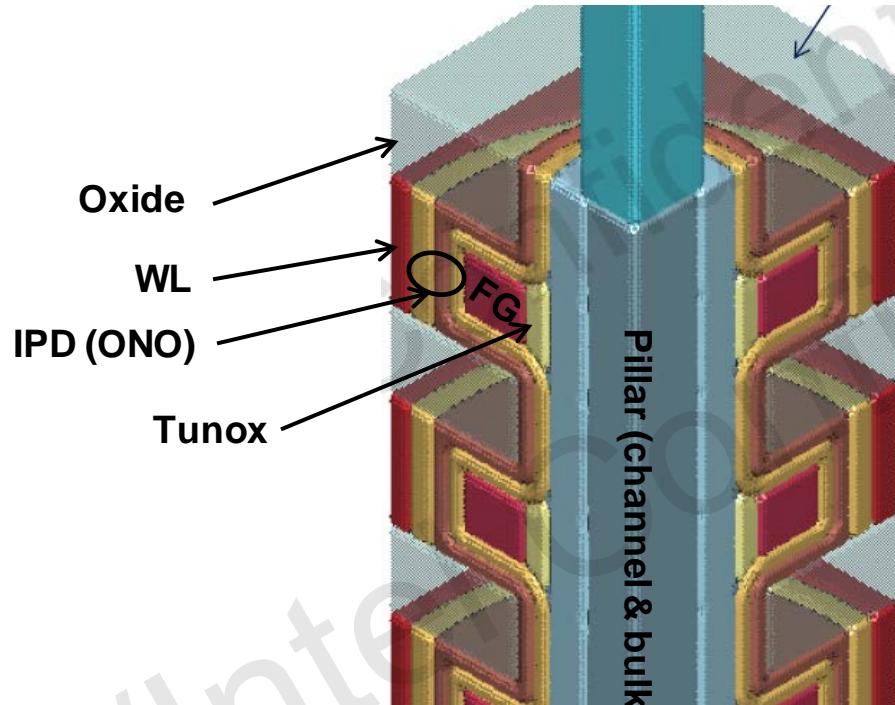
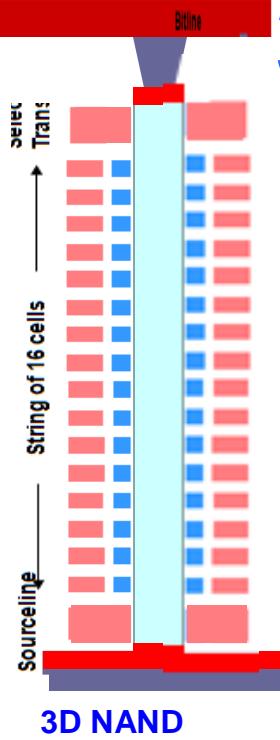
Only 16 cells/string shown for clarity.
80s (L85) has 128 cell string.



Only 16 cells/string shown for clarity.
Actual number of tiers is more.

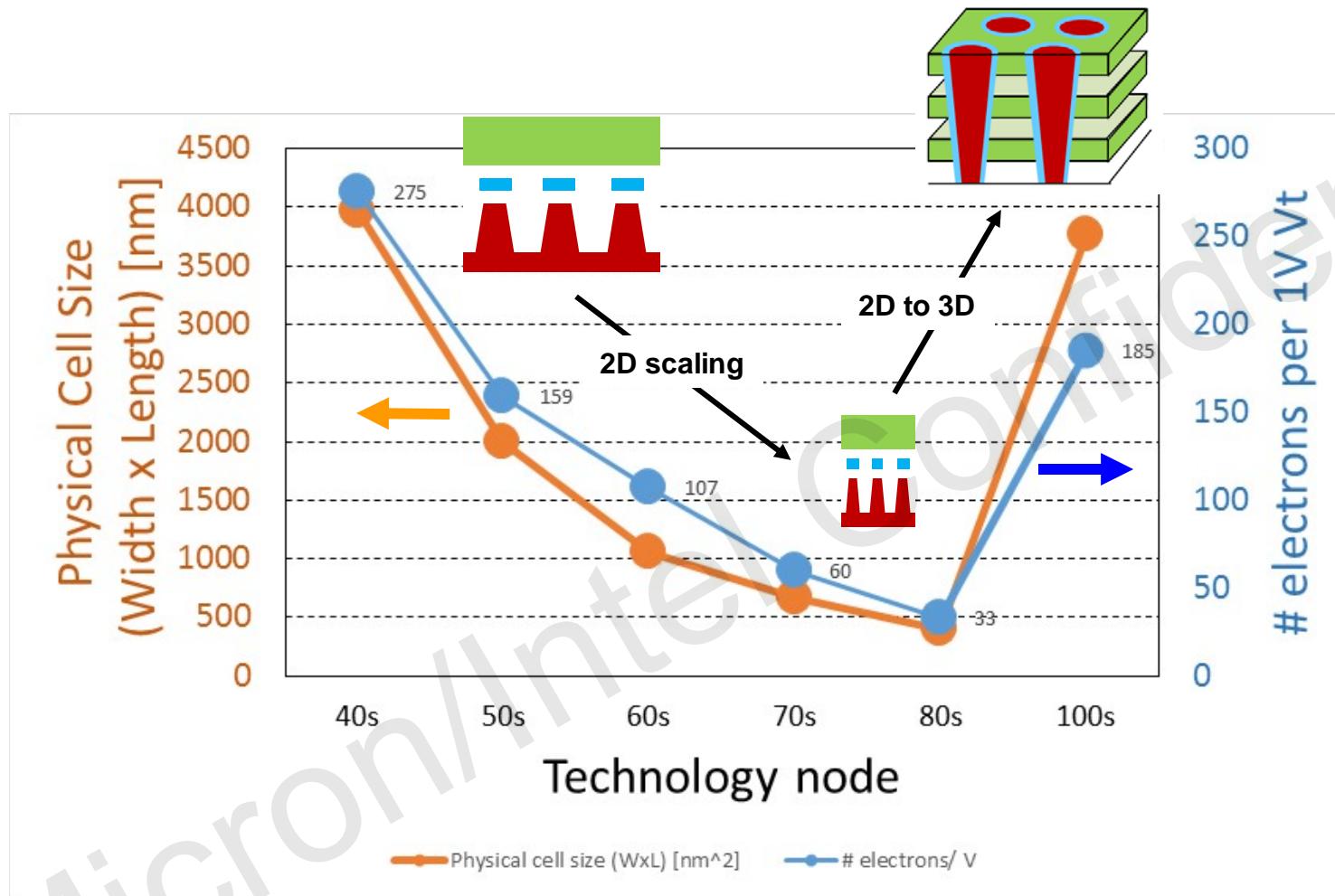
- Vertical NAND string
- Conductive or dielectric charge storage node
- Deposited poly-silicon channel
- New WL contacting scheme
- Large foot-print of the 3D Cell (~16X) will require quite a few layers to be stacked to achieve effective cell area scaling
- Increased demands on process technology – very high aspect ratio etches and fills

3D NAND Cell and String



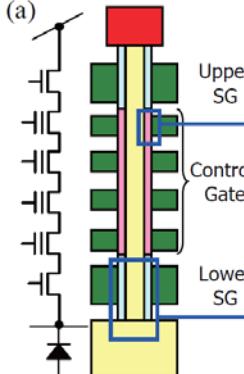
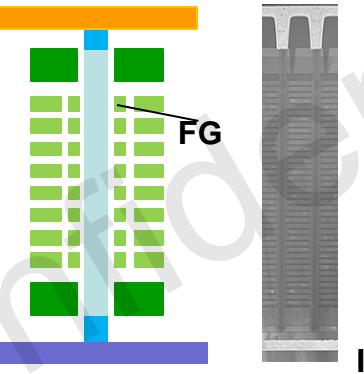
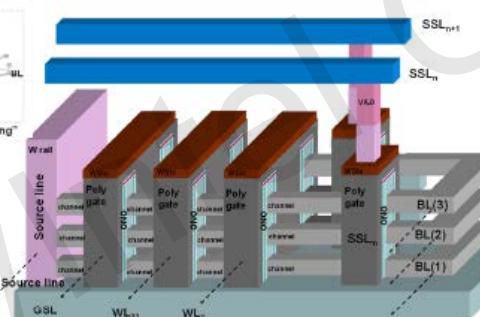
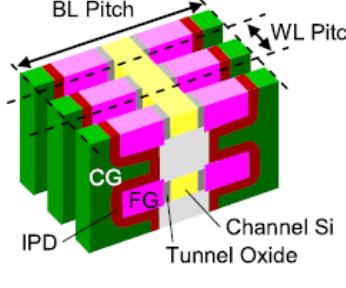
- Center of string is a cylinder of polysilicon: the “pillar”
- The pillar serves as the silicon bulk, and contains the channel
- The cell stack (Tunnel-oxide, Charge Storage Node, IGD, Control Gate) wraps around the pillar
- Adequate Gate coupling achieved through higher IGD/CG area

2D NAND to 3D NAND ~ Cell size~



3D NAND has a very large cell size owing to the gate all around architecture, realizing excellent electrical properties.

3D NAND Gallery: String and Cell

	SiN trap cell	Floating gate cell
Vertical string	 <p>(a) SiN trap H. Tanaka et al., VLSI 2007</p>	 <p>FG Intel-Micron</p>
Horizontal string	 <p>H.T. Lue et al., VLSI 2010</p>	 <p>BL Pitch WL Pitch CG IPD Channel Si Tunnel Oxide K. Sakuma et al., EDL 2007</p>

[Vertical String] Gate-All-Around, Cell properties advantage

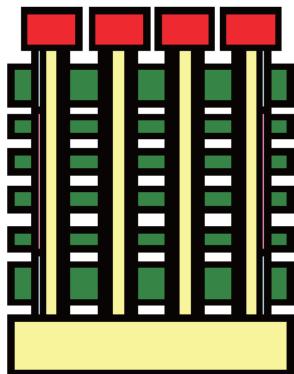
[Horizontal String] Small cell size, Tier stack height/cell foot print advantage

[SiN trap cell] Process integration & Low FGFG coupling advantage

[FG cell] Cell reliability advantage

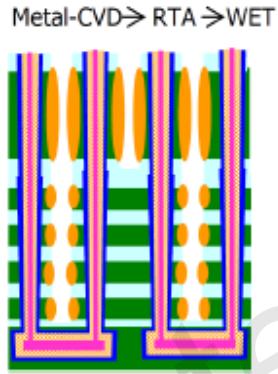
3D NAND Gallery: WL technology

Poly Si WL



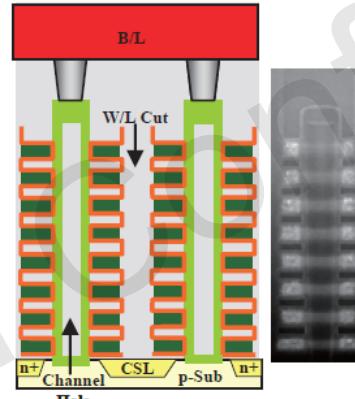
H. Tanaka et al.,
VLSI 2007

Salicided WL



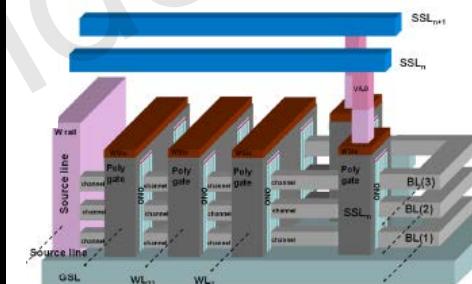
M. Ishiduki et al.,
IEDM 2009

Tungsten WL



J. Jang et al.,
VLSI 2009

WiSix WL



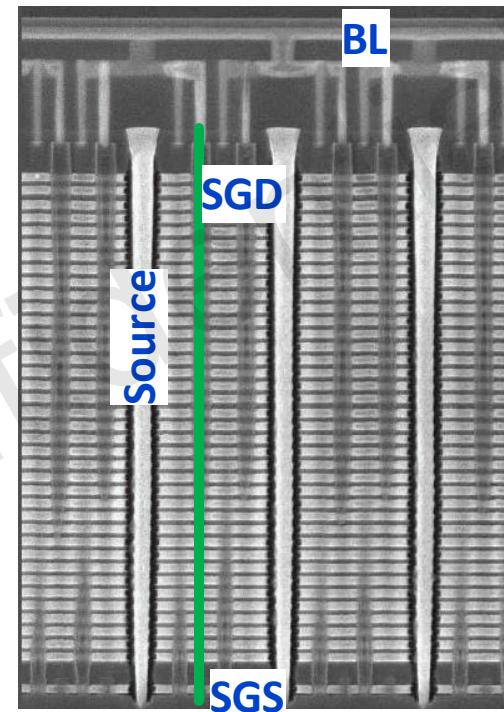
H.T. Lue et al.,
VLSI 2010

[Poly Si WL] Process integration and small foot print advantage.

[Low res. WL] Enabling long WL, string driver layout advantage.

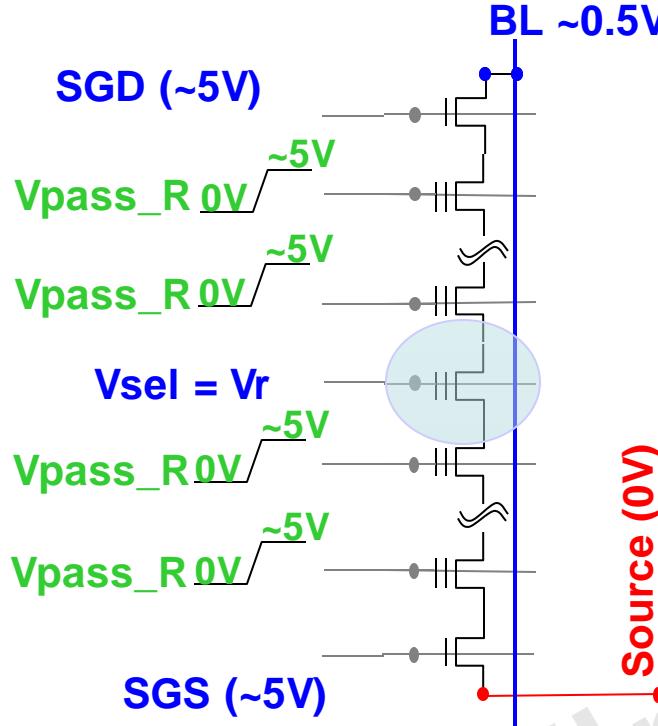
3D NAND Gallery ~ Production

NAND String

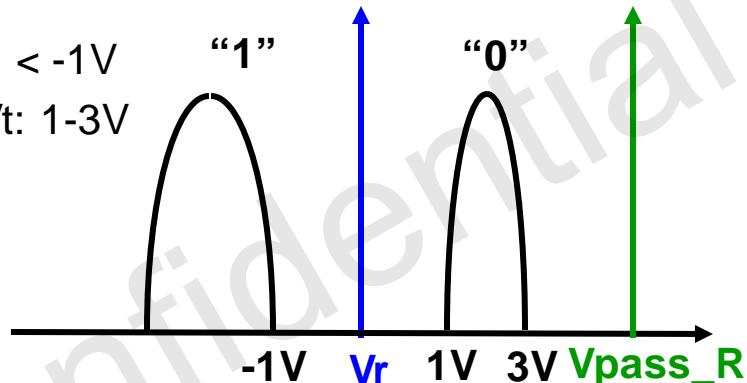


	Intel-Micron 3D FG	Hynix SMarT	Samsung/Toshiba V-NAND
Cell	FG	SiN trap (BE-TANOS)	SiN trap (BE-TANOS)
WL	Poly	Replacement W	Replacement W
SGD	MOS	Cell	Cell
SGS	MOS	Cell	MOS
String	Vertical, Straight	Vertical, U-shape	Vertical, Straight
CMOS	Under Array	Outside Array	Outside Array

NAND Flash Read



Erased Cell V_t : $< -1V$
Programmed Cell V_t : $1-3V$



- Bitline is biased to small voltage (~0.4-0.5V). SGD & SGS of the selected block (string) are turned on ($V_g \sim 5V$).
- Selected wordline is biased to Read Gate voltage (V_r). All the other wordlines in the string are biased to V_{pass_R} voltage (> the highest program V_t).
- If the selected cell $V_t < V_r$, the string will conduct and there will be current on the bitline. If the selected cell $V_t > V_r$, the string will not conduct and the current on the bitline will be zero. Sensing circuits detect this.

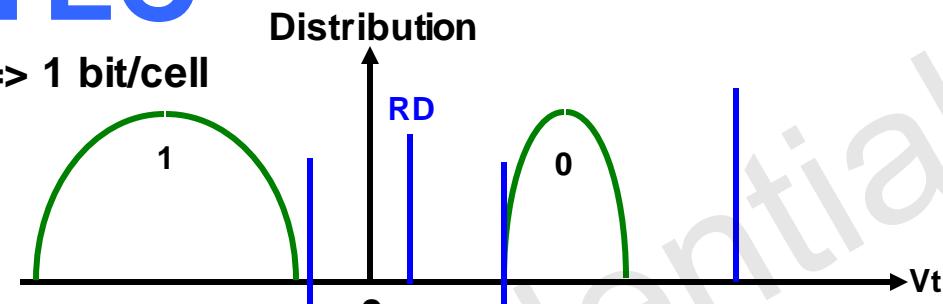
SLC, MLC, & TLC

SLC: 1 bit/cell: 2 states:

Level 0 = Erased

Level 1 = Programmed

2 Levels => 1 bit/cell



MLC: 2 bits/cell: 4 states:

Level 0 = Erased

Level 1 = Pgm to V_{t1}

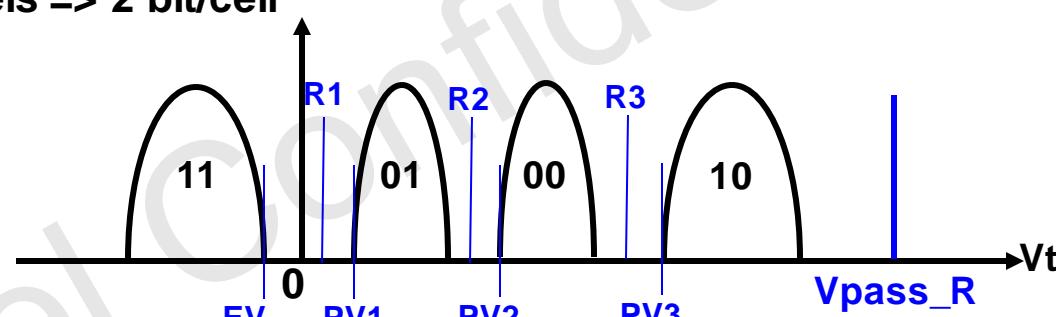
Level 2 = Pgm to V_{t2}

Level 3 = Pgm to V_{t3}

LSB (LP) Read with V_g = R₂

MSB (UP) Read with V_g = R₁ & R₃

4 Levels => 2 bit/cell



TLC: 3 bits/cell: 8 states:

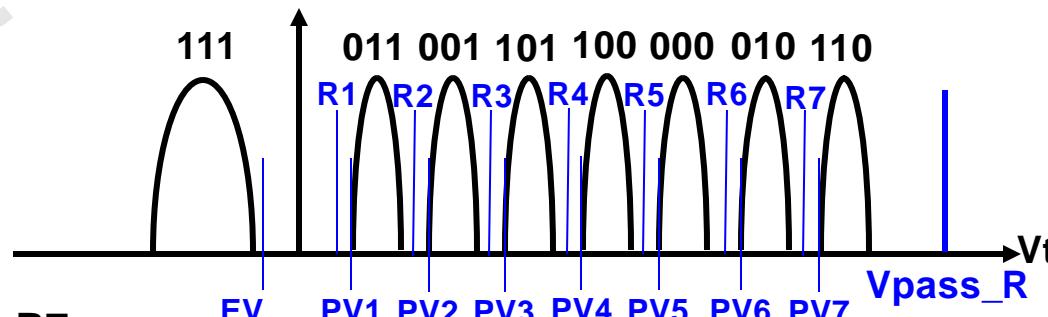
Level 0 = Erased

Level 1-7 = Pgm to V_{t1} – V_{t7}

LSB (LP) Read: Read with V_g = R₄

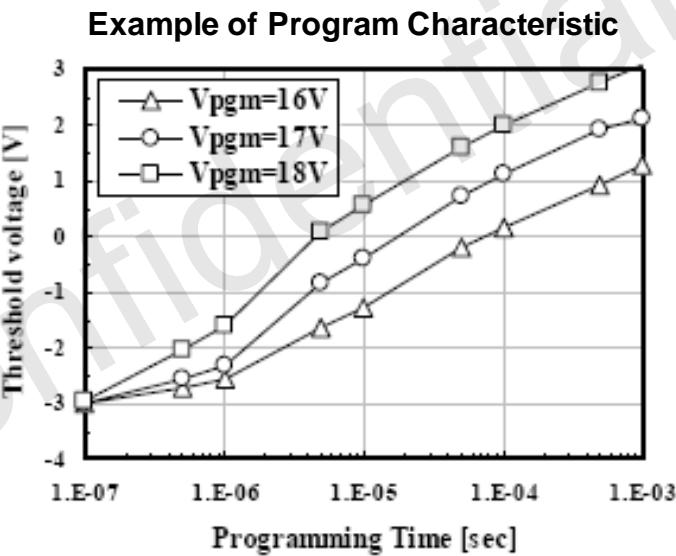
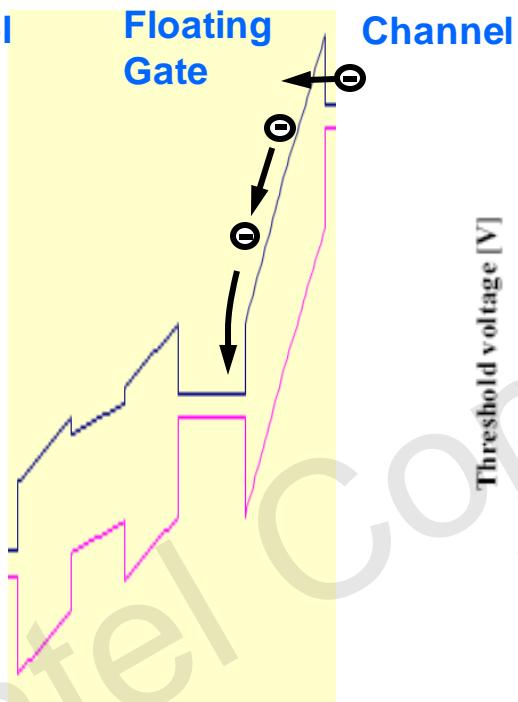
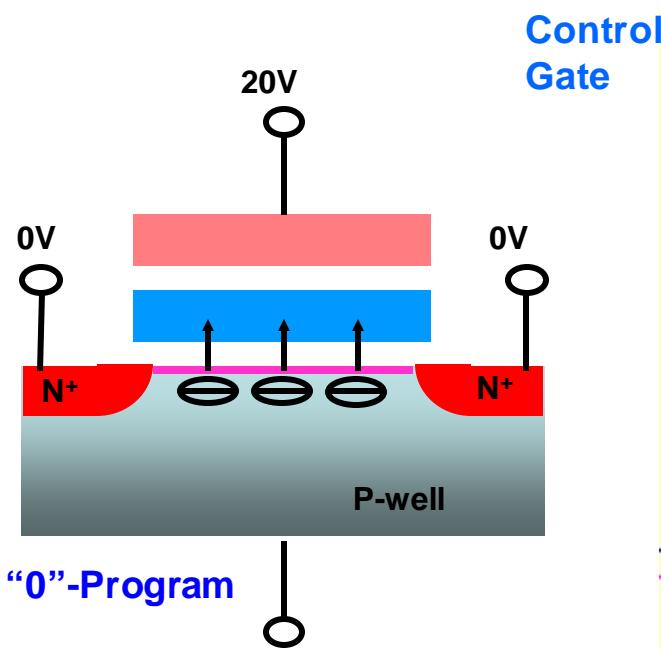
?SB (UP) Read with V_g = R₂ & R₆

MSB (XP) Read with V_g = R₁, R₃, R₅, & R₇



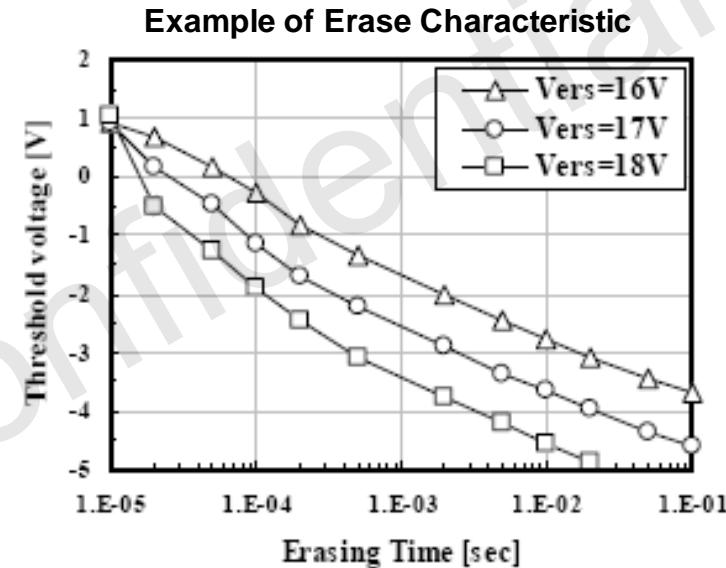
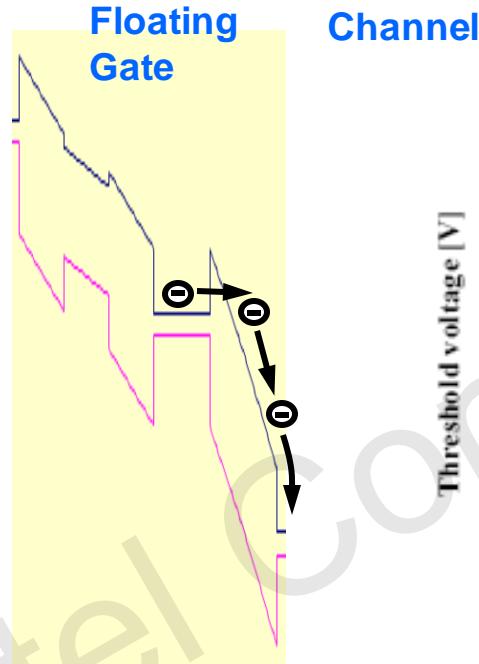
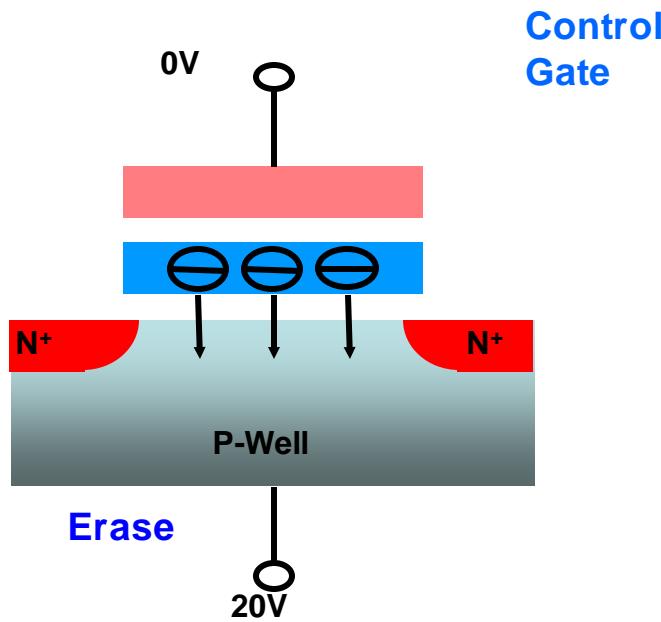
- Having the threshold voltage of the cells defined to 2^n different levels allows for storing n bits per cells

NAND Flash Programming - FN Tunneling



- Programming is by Tunneling in of electrons from the channel into the FG by biasing the Control Gate positive w.r.t the channel and S/D.
- “Effective” Program Time (cumulative pulse time ignoring set up and verify Overheads) ~100-200us.
- Tunnel-oxide E-field ~11-12MV/cm. Tunneling Current ~ 10mA/cm².
- Program page size ~ 16KB

NAND Flash Erase - FN Tunneling

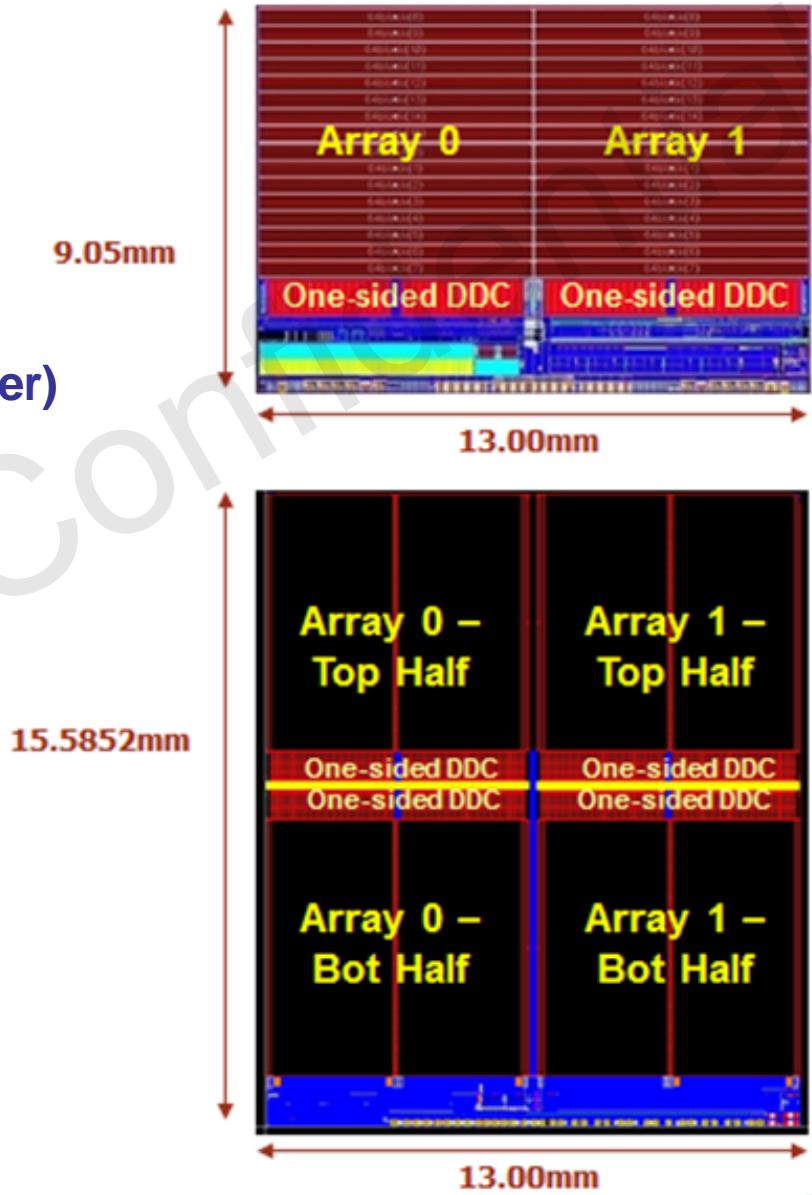


- Erase is opposite of Programming by reversing the applied biases.
- Tunneling out of electrons from the FG into the channel by biasing the Well positive w.r.t the Control Gate. Erase time ~1-10ms.
- Tunnel-oxide E-field ~11-12MV/cm. Tunneling Current ~ 1-10mA/cm².
- Erase Block size ~ 24MB/96MB/36MB (B0KB/B1MA/B16)

More on Read, Program and Erase Later

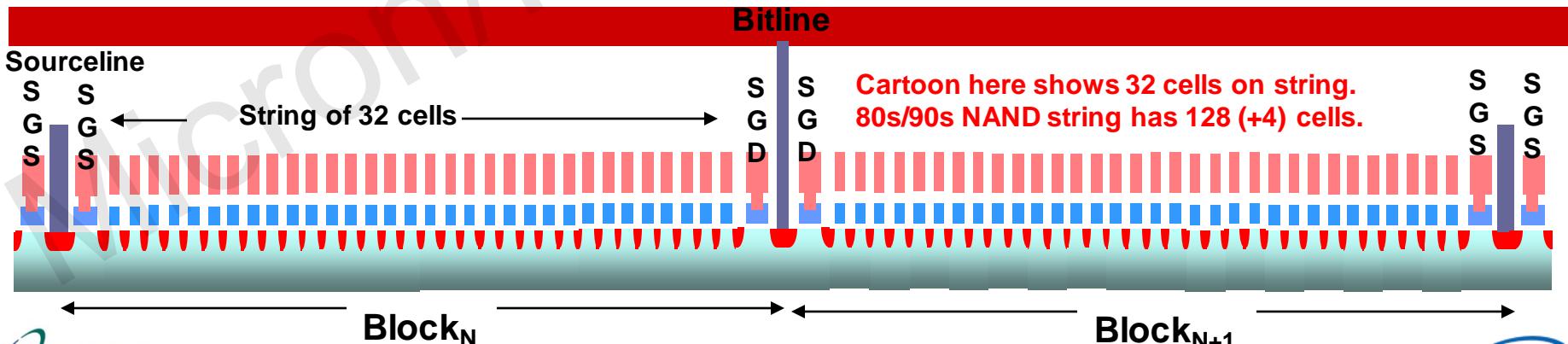
2D NAND Array – 80s

- L84A
 - ▶ Dual Plane
 - ▶ Bottom 1-sided 8KB DDC
 - ▶ 68 NAND String
 - ▶ Mixed String Driver (Edge & Center)
 - ▶ ECC: 1KB code word
 - ▶ Pad arrangement: 1-sided long
- L85A
 - ▶ Dual Plane – Folded Array
 - ▶ 1-sided 8KB DDC for half plane
 - ▶ 132 NAND String
 - ▶ Staggered Mixed String Driver
 - ▶ ECC: 2KB code word
 - ▶ Pad arrangement: 1-sided short



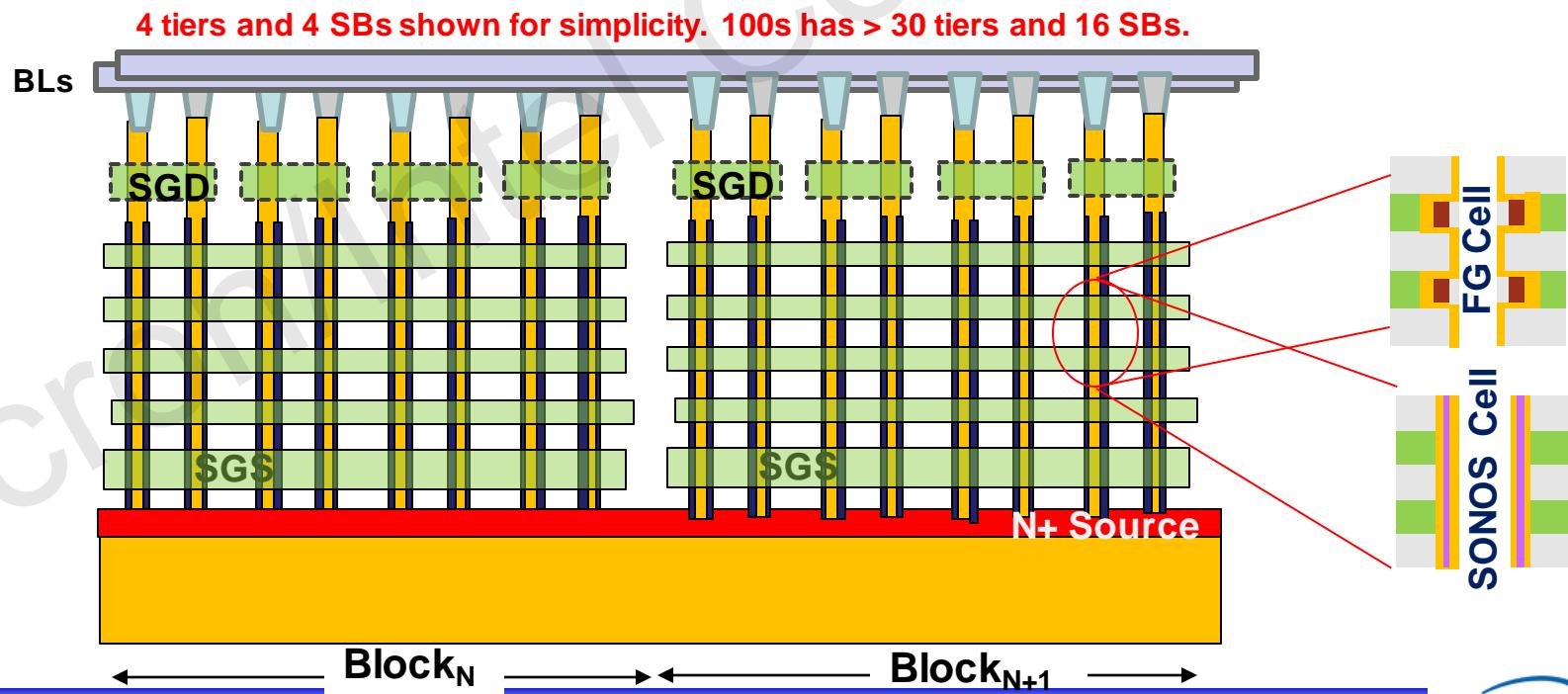
2D NAND Array

- Die is typically divided into 2 or 4 planes
- All blocks in a plane typically sit in a single P-Well in a single Deep N-Well (Referred to as ATUB)
- Bitlines (columns) typically span the full length (or half length) of the die
- Wordlines typically span the full width (or half the width) of the die
- Blocks are separated by the select transistors. A given string or a block is selected by selecting the appropriate Select Gates (SGS & SGD)
- Bitlines contact the drain side of the Select Gate Drain
- Source is common to the entire plane
- A block has 128 Rows (Wordlines). With 64-128 cells in the string, one Block = $(64-128 \text{ WLs}) * \# \text{ of Cell on the WL}$ (8-16KB)

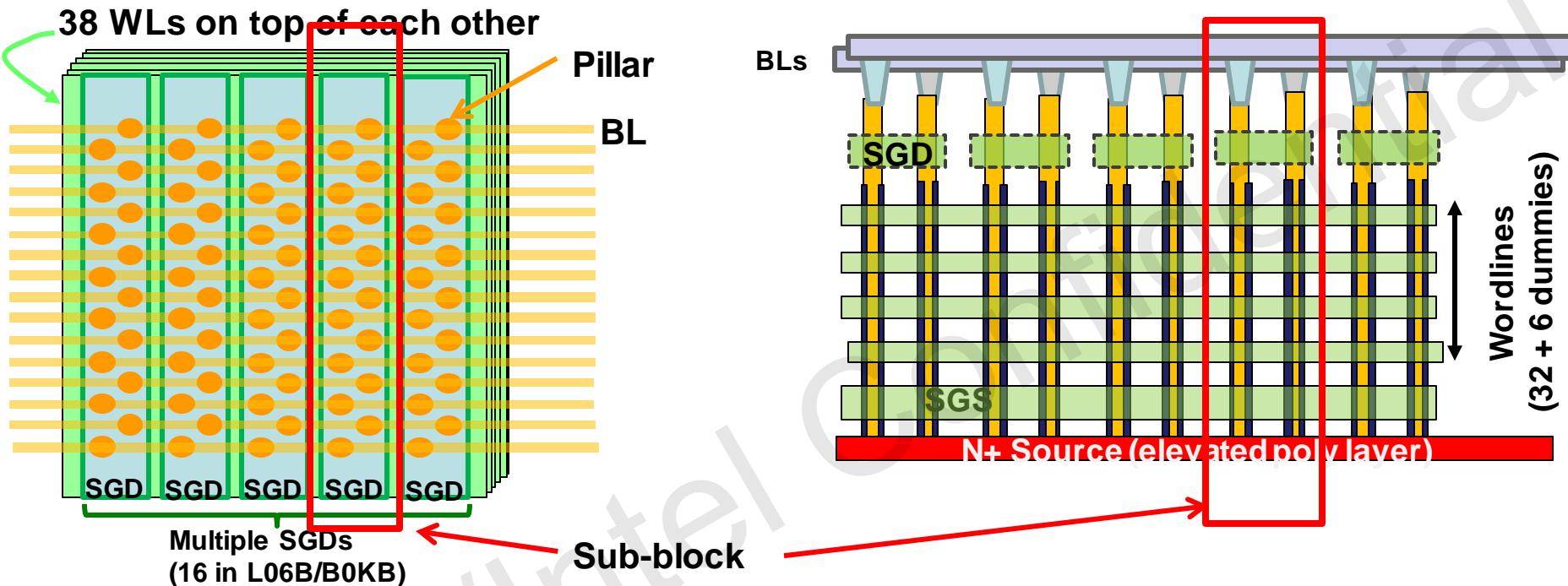


3D NAND Array (L06B/B0KB)

- 3D NAND introduces additional unit of array – sub-blocks. An erase block is divided into 16 sub-blocks (L06B/B0KB) – selected by the SGD.
- Wordlines are shared across the sub-blocks in a block.
- SGS is shared across the sub-blocks in a block.
- Pillar is the NAND string and the Atub.
- L06B/B0KB: With 32 tiers, 16 sub-blocks (SGDs/block), and 16KB page size, one Block = (32 WLs) * # of Cell on the WL (16KB) * 16 sub-blocks



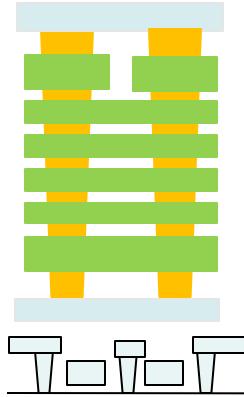
3D NAND Block (L06B/B0KB)



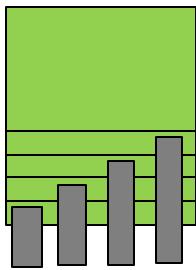
- Wordline: Many cells wide (16 in L06B/B0KB)
 - Ideally would have WL just one cell wide (as in 2D), but cannot fit string drivers/WL hook-ups for all the WLs into that small a pitch
 - So WL is wide to fit the string drivers and individual NAND strings are selected by SGDs – Multiples SGDs in a block → Many Sub-blocks / Block
- One page = intersection of one WL and SGD
- One cell = intersection of one WL, one BL, and one SGD

What Decides SGDs/Block

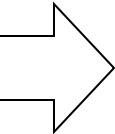
CMOS under Array
(Cross sectional view)



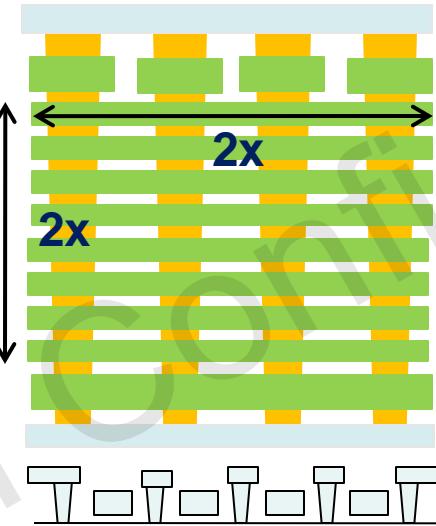
WL exit routing
(Plane view)



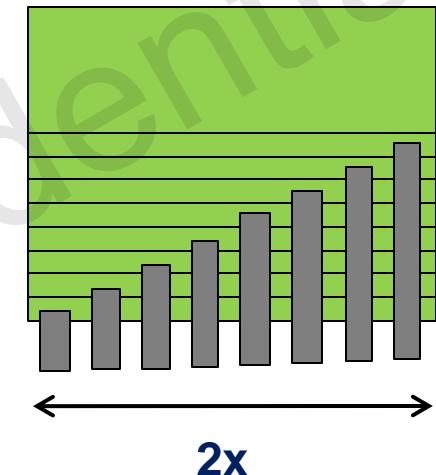
2x increase
of stack



CMOS under Array
(Cross sectional view)



WL exit routing
(Plane view)



- Cell is 3D while CMOS circuit and routing are 2D.
- Cell 3D scaling (stacking) requires more 2D area for circuit and routing.
- This results in # sub-block(SGDs) increase → increased disturb stress and WL capacitance. WL RC remains unchanged.

3D NAND Pillar/SGD Layout

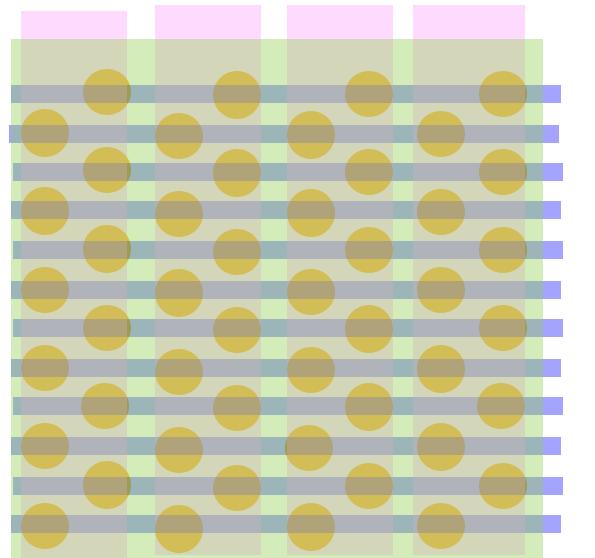
Simple Case

- Simple arrangement of pillars in square grid.
- BL Pitch = Pillar Pitch
- SGD Pitch = Pillar Pitch
- Pillars are not closely packed → Poor array efficiency
- Smaller Page size → Lower Program Bandwidth
- Block Size (WL width) decided by the String Driver layout needed to support WL and SG drivers for all the tiers and SGs

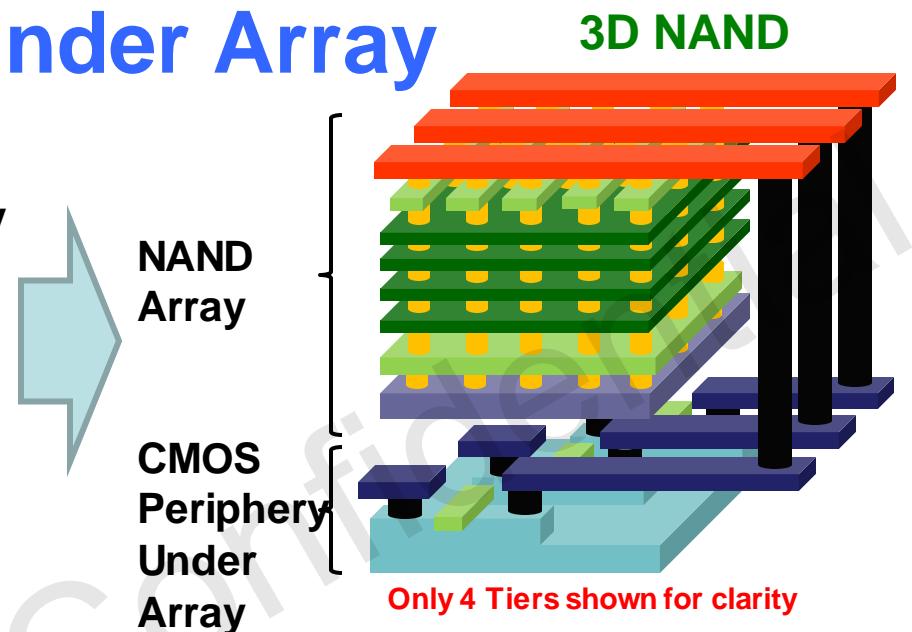
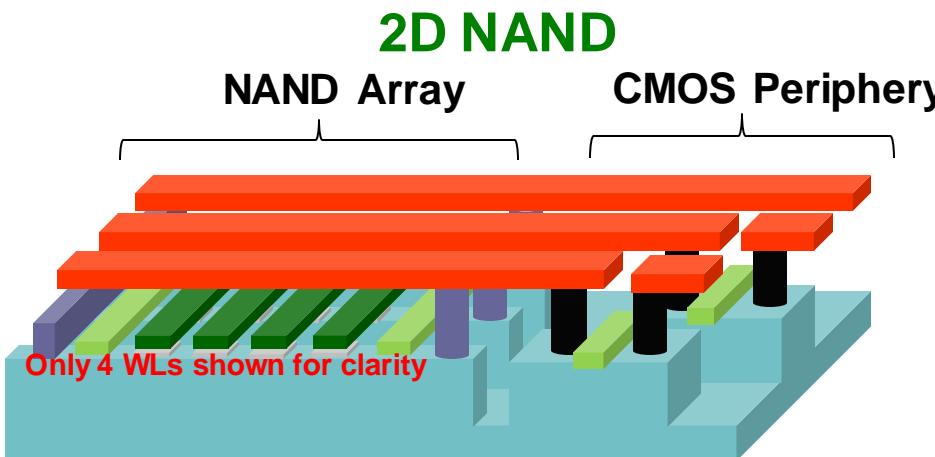


Stagger the Pillars (L06B/B0KB and L17A/B1MA):

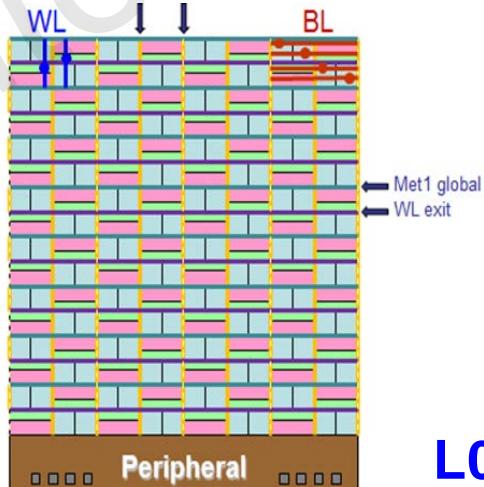
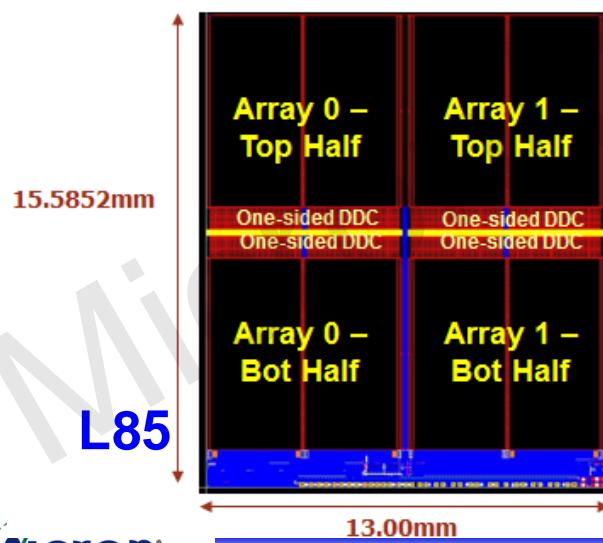
- Hexagonal Close Packed Pillars → Best array efficiency
- BL Pitch = $\frac{1}{2}$ Pillar Pitch → Capable of 2X page size → Improves the Program bandwidth by 2X
- ~Doubles the SGD Pitch → For a given WL width (dictated by String Driver layout) number of SGDs/sub-blocks is halved.



3D NAND with CMOS Under Array

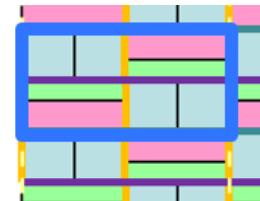


- 3D NAND is very compatible with CMOS/Periphery Under Array → Sense Amps, Array drivers placed under Array → Tile Architecture.

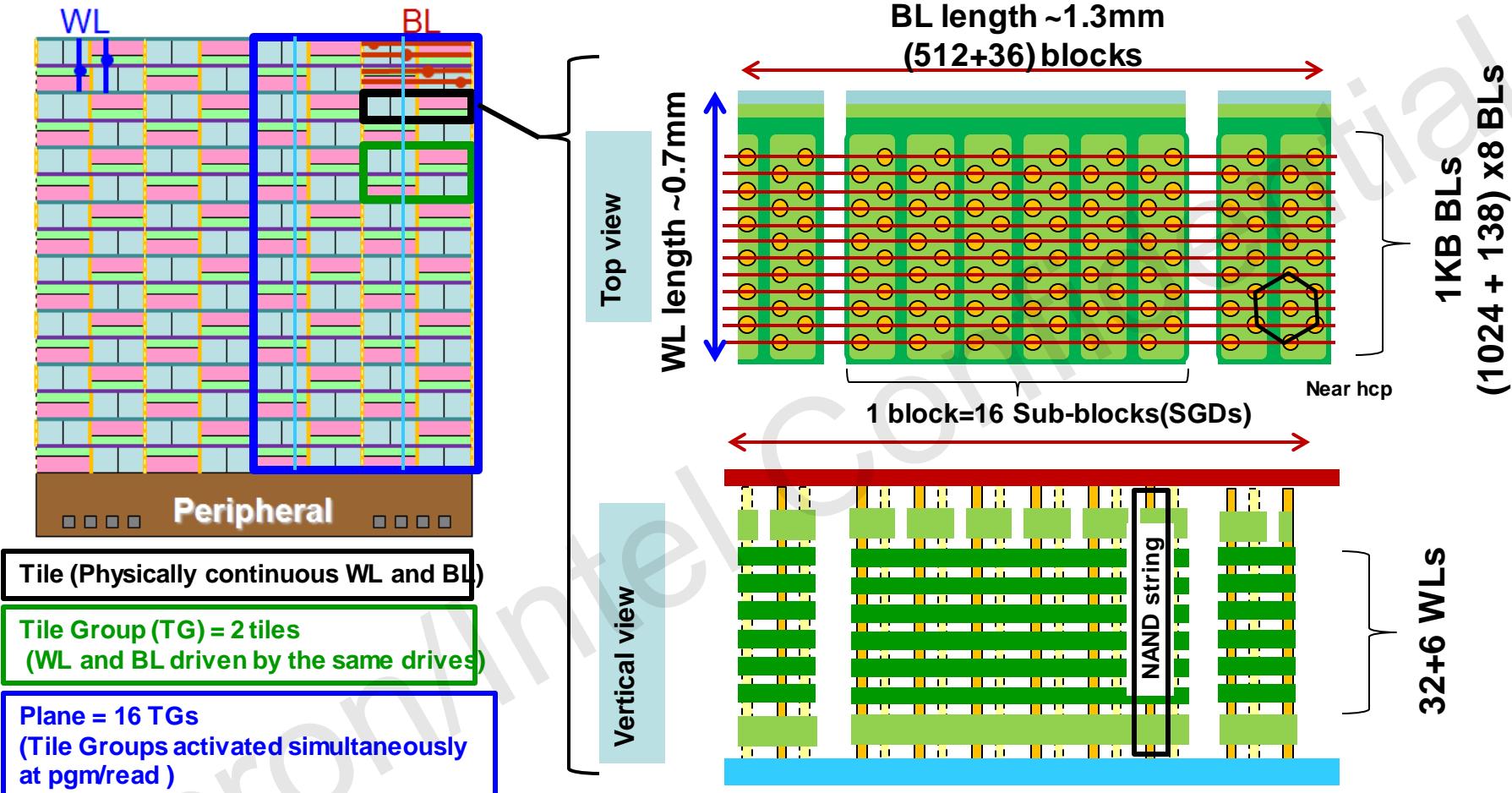


1 Tile Group (TG)
19008 BL x 548 block segments
Spans 4 SRC plates
are separated by Staircase/WL Exit
there are 32 tile groups per die (8 x 4)

L06A



3D NAND Architecture



Page size: Number of cells program or read together

= number of BLs in a plane

= 32KB (1KB/tile x 2 tiles/TG x 16TGS/plane)

Pages per block: 1024 pages (=2pages/WL x 32WLs/sub-block x 16 sub-blocks/block)

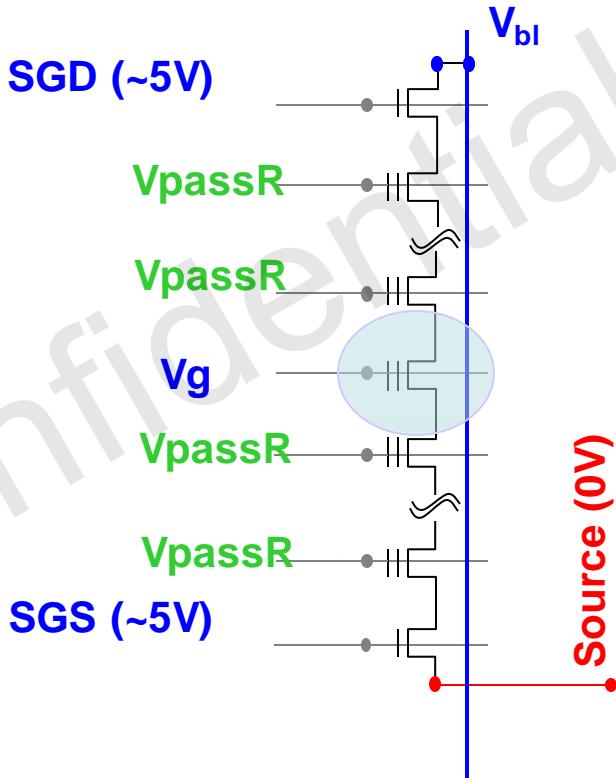
Die density: 256Gb (=32KB/page x 1024pages/block x 512blocks/plane x 2planes/die)

Read

- Cell/String I-V
- String Current Requirement
- 100s String Current limiters
- No-LDD issues
- Sensing – SBL & ABL
- Interference (FG-FG)
- 1-Pass and 2-Pass Programming
- Shifted Window

NAND Cell I-V

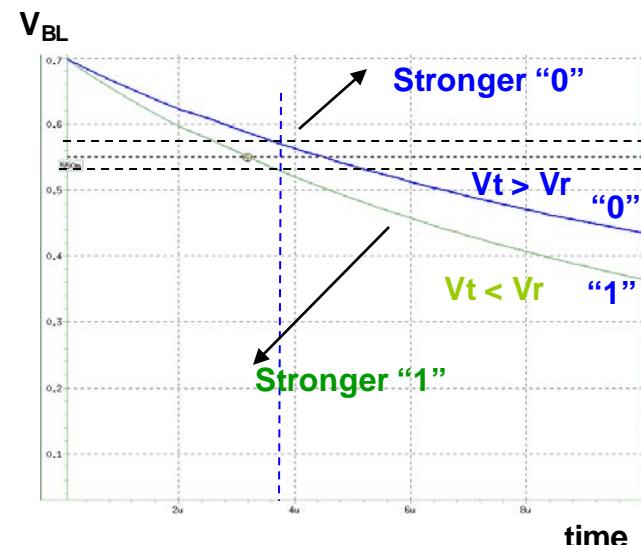
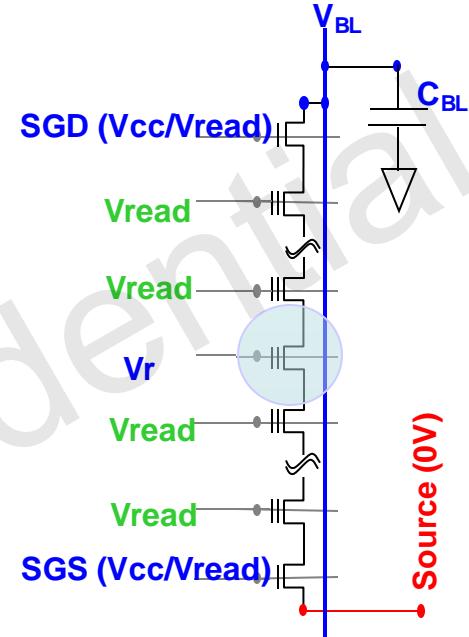
- A NAND String can comprise 32-128 NAND active cells plus two select transistors
- Like any transistor, there is gate bias (V_g) and Well/substrate bias (V_b) which can be directly applied for the selected cell
- However, there is no direct access to the source and drain of the selected cell. External bias can only be applied to the Bitline (V_{bl}) and the Sourceline (V_s)
- The select gates (SGS & SGD) need to be biased to a high V_g (~5V) to turn these on
- The deselected cells need to be turned on as well by making sure that these gates are biased well above the V_t of these cells
 - The deselected cells could have any V_t (Back-Pattern: BP) depending on whether they are in the Programmed or Erased state. The Cell V_t s and the gate biases will dictate their channel resistance. This in turn will decide the actual source and drain bias the selected cell sees.



Basic Sensing Concept

Simple description of voltage-sensing concept (pre-100s)

- Selected bitlines are charged to a V_{BL_init} and then floated
 - String to be read is selected:
 - Wordline to be read is biased to the Read Gate bias (V_r) and all the other wordlines in the string are biased to V_{pass_r} to make sure they are able to conduct
 - SGD and SGS are turned on.
 - The cell current will start discharging the bitline
 - After a certain amount of time ($t_{SAD} \sim 5-10\mu s$) the Bitline voltage is sensed.
 - If the $V_{BL} > V_{BL_fin}$ (not discharged) $\rightarrow "0"$
 - If the $V_{BL} < V_{BL_fin}$ (discharged) $\rightarrow "1"$
 - Let $\Delta V_{BL} = V_{BL_init} - V_{BL_fin}$
 - Then, if $I_{cell} > \Delta V_{BL} * C_{BL} / t_{SAD} \rightarrow "1"$
 $I_{cell} < \Delta V_{BL} * C_{BL} / t_{SAD} \rightarrow "0"$
 - For selected cell to be sensed, current should be dictated by selected cell and not rest of string.
- The string current $>>$ sense current.



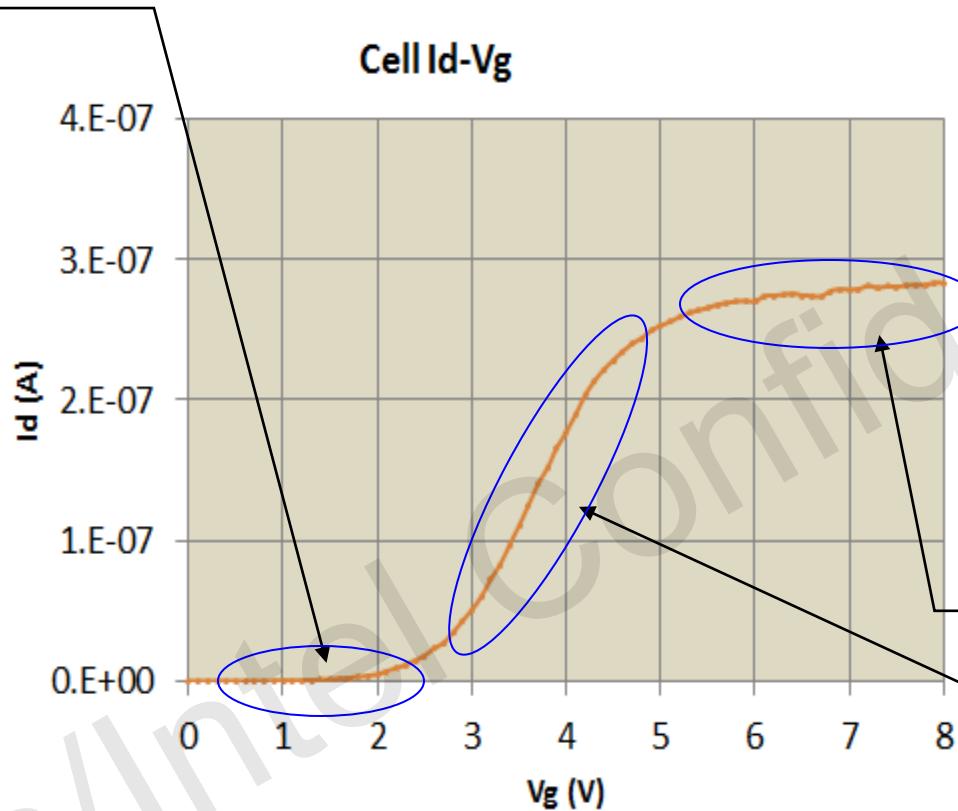
String Current Required

#	Item	Value	Consideration
1	Sense current	~10-15nA	Margin to leakage, Noise, Sense times and BL Stabilization (Performance).
2	Dual Strobe Margin	~1.5X-2.0X	PPV sense current is 1.5-2X higher
3	Sense current to String current Guardband	~3X	Current set by Cell Vt and not string, since we want to read cell. RWB.
4	Temp Guardband (90C to -10C)	~1.5X	Poly channel current degrades at low temperature (in single crystalline silicon current increases at low temp)
5	Cycling Guardband	~1.5X	Trap-up related current degradation
6	Erase string to Programmed String	~1.7X	Reduced Gate overdrive on Programmed string vs Erased string
7	90C Probe Min Erase string current Required	~120nA - 200nA	

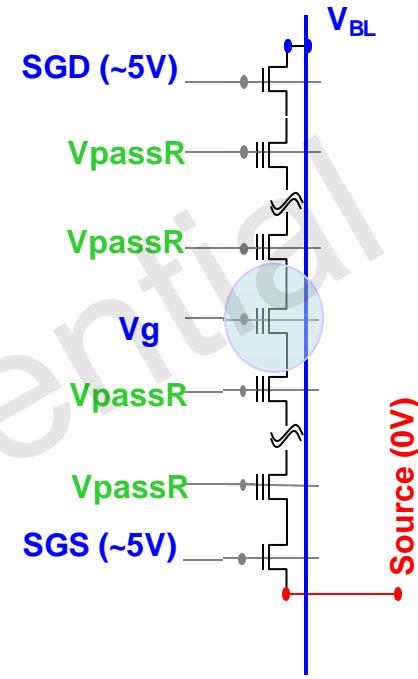
- Intrinsic Cell improvements (Temp/Cycling GB) and Process improvements (DHC interfaces, distributions, leakages) can improve many of these factors to manage the overall string current requirement.

NAND Cell I-V

Current limited by the cell. Full Bitline to Source voltage drop occurs across the selected cell.
Assuming the Bitline voltage is V_{BL} and the source is at 0V, the selected cell will see $V_{ds} \sim V_{BL}$.

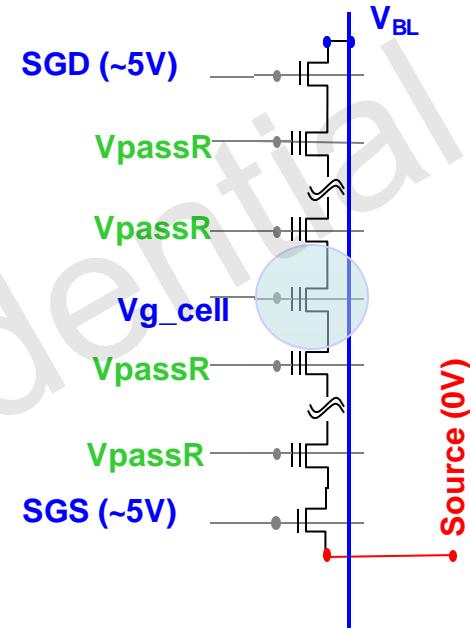
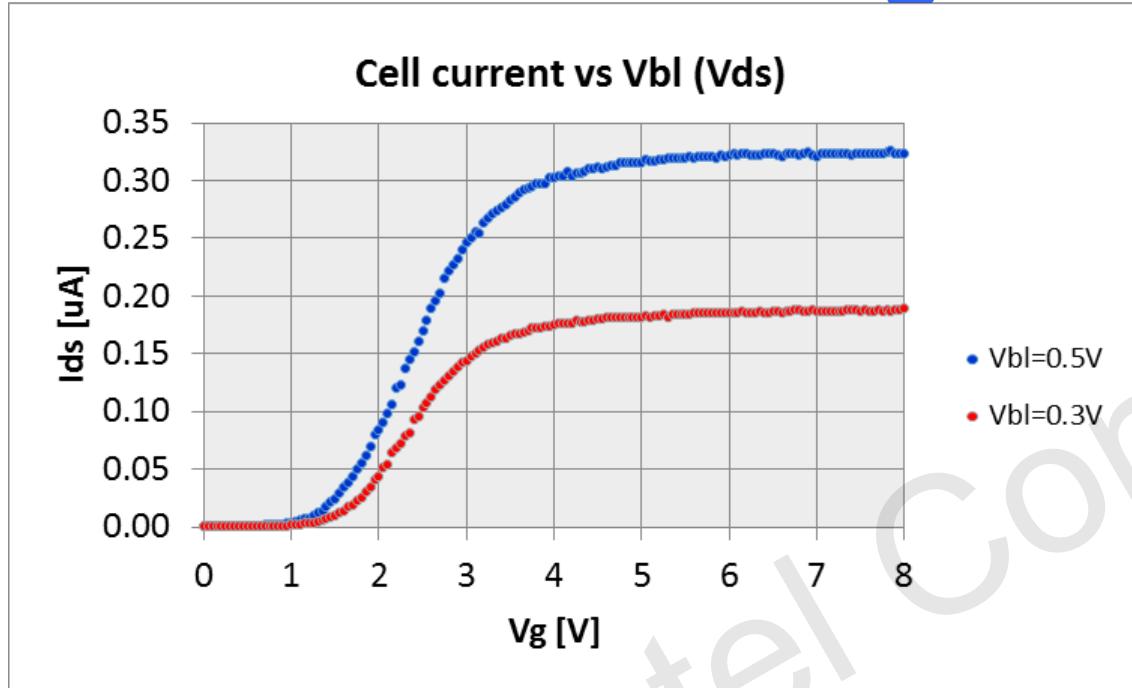


As the cell starts to conduct, there will be IR drop in the rest of the string causing the V_{ds} across the selected cell to decrease. If the cell is on a lower wordline IR drop is mostly on the drain side causing the V_d to decrease. If the cell is on a higher wordline, the IR drop is mostly on the source side causing the V_s to increase. Since the V_s , V_{ds} , V_{gs} for the cell are all varying the slope of the I-V does not really represent the cell G_m .

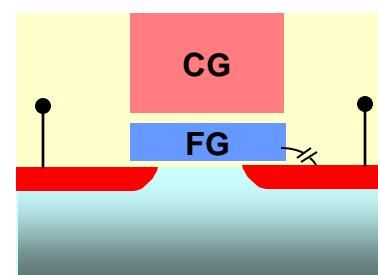


Saturated string current. Current limited by the total string resistance ~ LDD resistance as well as the net channel resistance of all the cells in the string and select gates. The selected Cell only sees a fraction of V_{BL} .

Cell I-V: BL Voltage

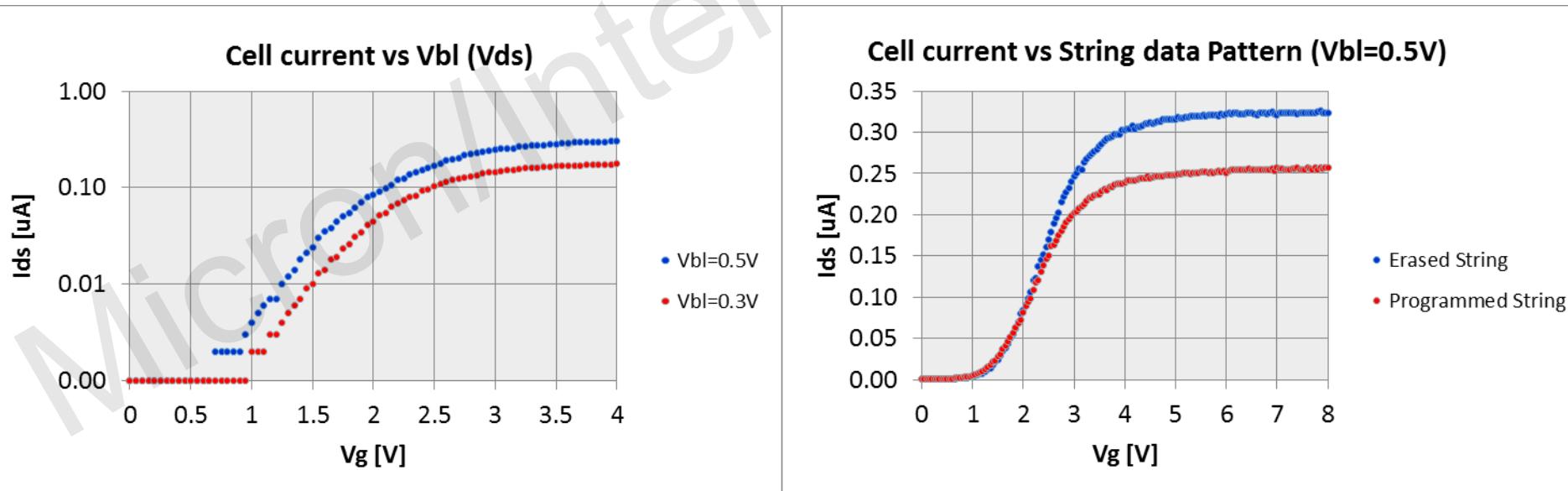


- Bitline voltage impacts the string saturation current as the string current is Bitline voltage/Net String resistance (Channel + S/D resistance of all the cells and SGs on the string).
- It also causes the Cell V_t to shift due to DIBL/DCR (Drain-induced Barrier Lowering / Drain Coupling Ratio).



Key Cell I-V Parameters

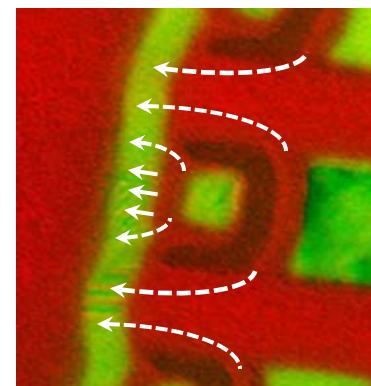
- Some of the key cell parameters of interest are:
 - Sub-threshold slope – can be measured from the I-V
 - DCR/DIBL: Change in the cell V_t due to change in the V_{ds}
 - Body effect: Change in the cell V_t due to change in the V_{sb} . On 3D NAND with floating body, there will be no body-effect. But change in V_s will impact $V_{gs}-V_t$.
 - Trans-conductance (G_m) – difficult to measure directly since the cell V_d & V_s in a string dynamically changes causing the V_t , $V_{gs}-V_t$ to change as well
- Knowing the above for a cell allows us to construct the cell/string I-V



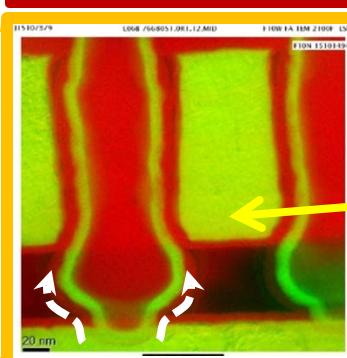
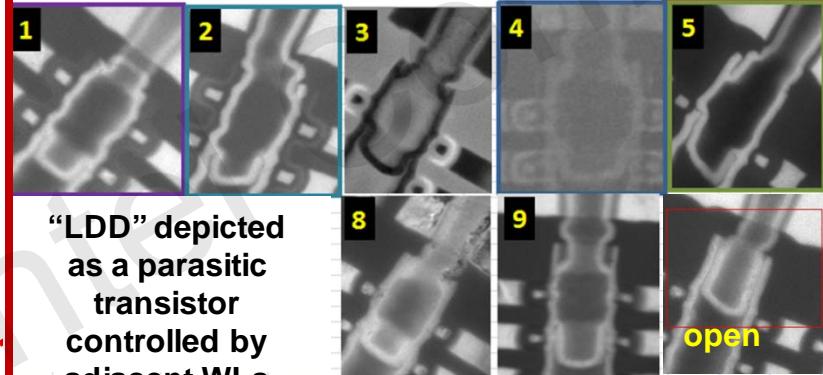
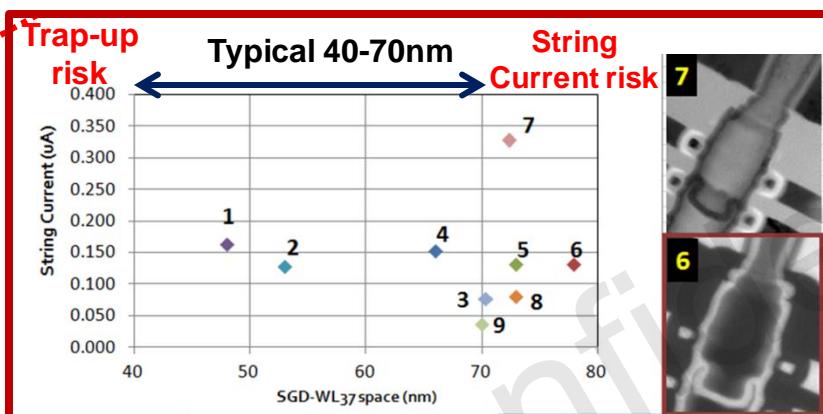
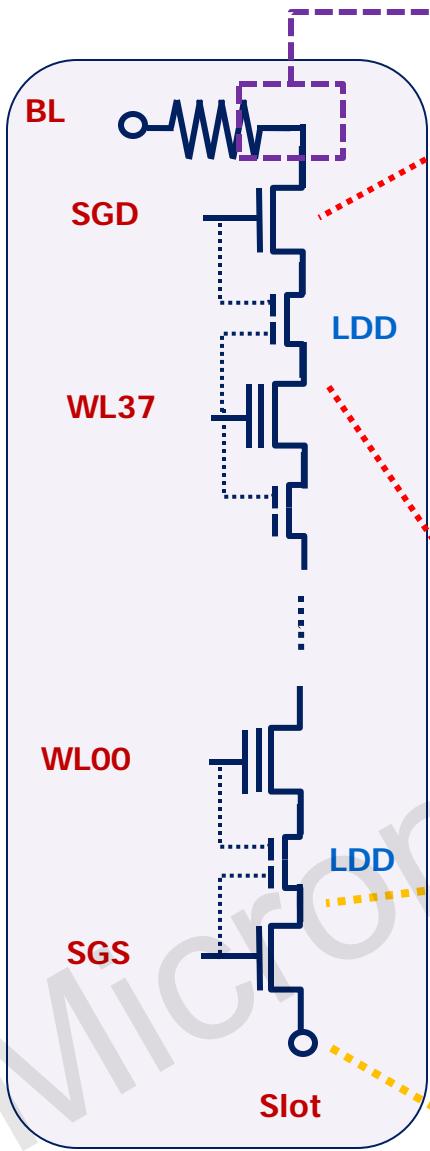
100s String Current

- String current health is one of the key parameters. There are three important factors to consider:

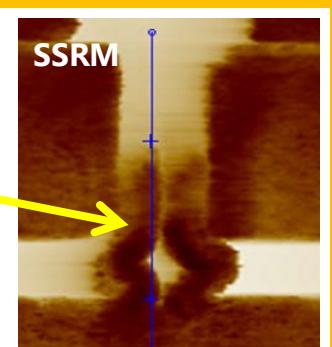
1. Channel material is made from poly-silicon. Inter-grain boundary and associated traps influence charge mobility significantly (~20-40x lower mobility compared to crystalline silicon). Moreover, the channel material is made of thin and hollow cylinder (implemented to improve cell V_t uniformity); charge transport is influenced by proximity of two oxide-poly boundaries
2. Cell transistors in the string are junction-less; the classic “LDD” region is not doped and relies on fringing field to operate (a less efficient coupling than the channel below cell)
3. Process integration such as 21L-24L, 24-41L, 21-45L are critical



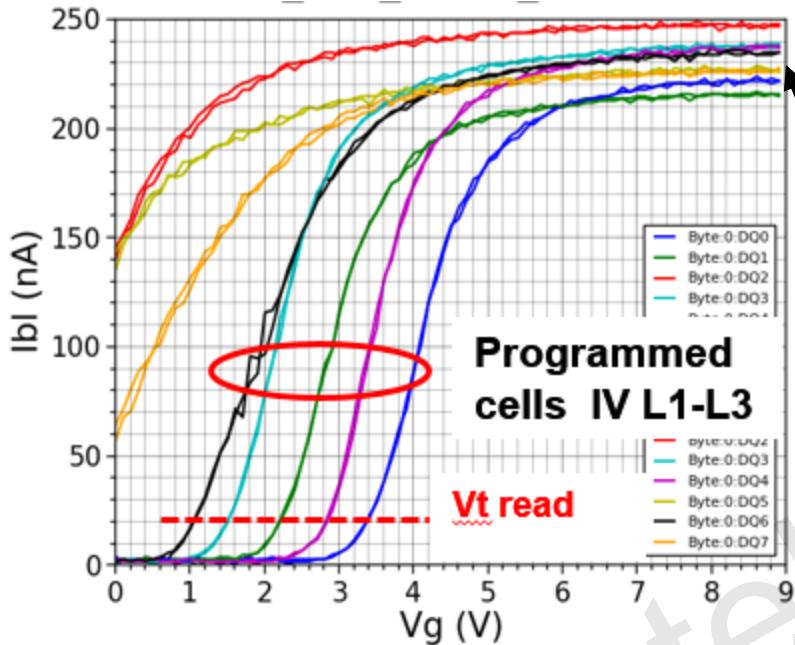
String Current: Process Integration



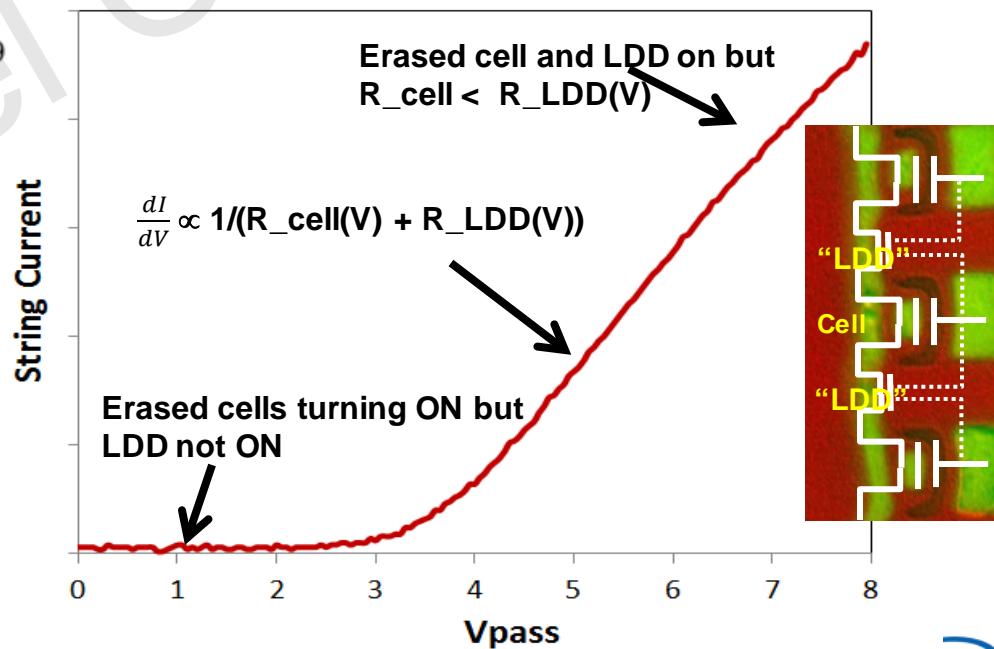
1. Dopant out diffusion from highly doped slot lowers LDD resistance. Too high doping lower SGS V_t (risk).
2. Main junction for erase GIDL generation.
3. Electron injection point and key for string current performance.



Cell and String I-V

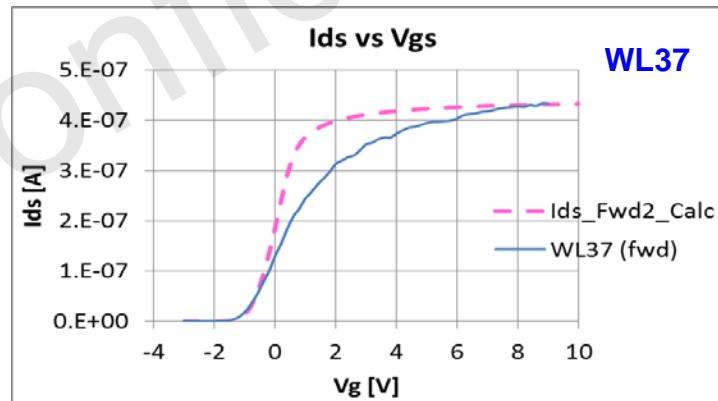
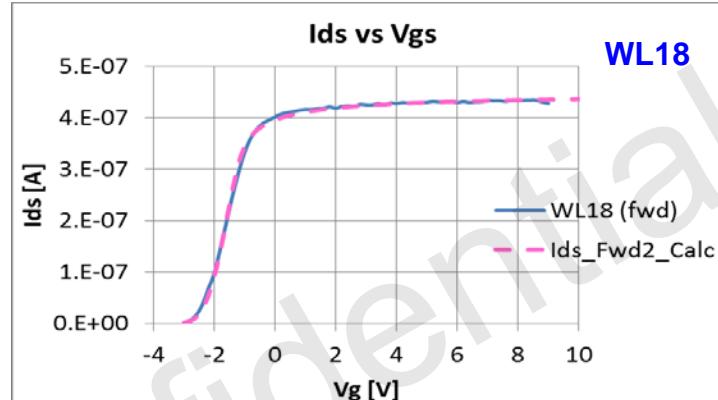
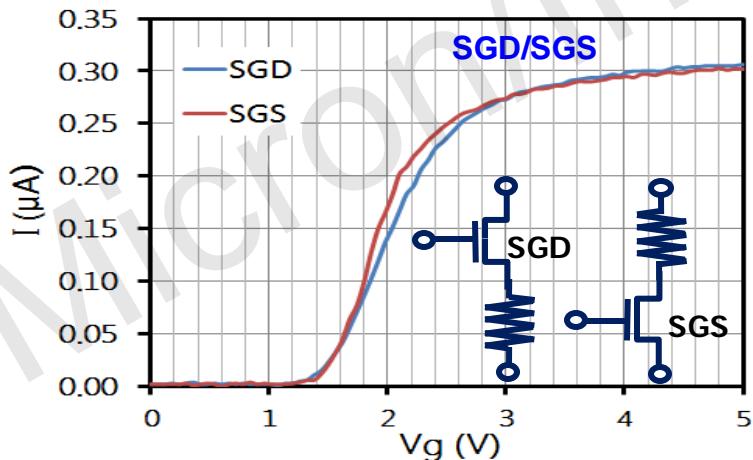


- Cell I-V (all WLs at $V_{pass_R} \sim 8-9V$, and sweep of selected WL) shows a steep turn-on of transistor (decided by Cell V_t) and saturation decided by string current.
- Limited by resistance “external” to the selected cell (Vts of other cells in the string, string to string variability, etc.)



Cell and SG I-V

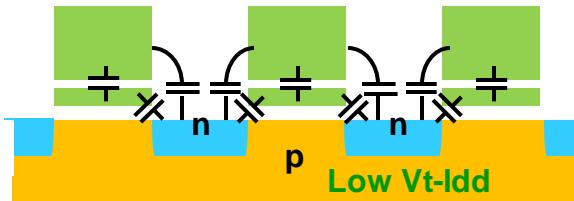
- Cell I-V can be modeled very well using simple transistor I-V equations (example WW18 Id-Vg shown here).
- Channel mobility is estimated to be ~5-10cm²/V-sec (compared to 100-200cm²/V-sec) for single crystalline Si (e.g. 80s).
- But the larger W/L ratio of 100s NAND string compensates for the reduction in mobility to provide adequate string current.
- WL37 example shows significant deviation of experimental I-V from expected I-V and → Softer saturation of current (significant voltage response) → indication that resistance between SGD-WL37 is still significant.



- Apparent “gm” of SGD transistor is lower than SGS as the main loading resistance is sitting on the source side of SGD transistor
- Loading resistance analysis can be helpful to segment location of string current limiters

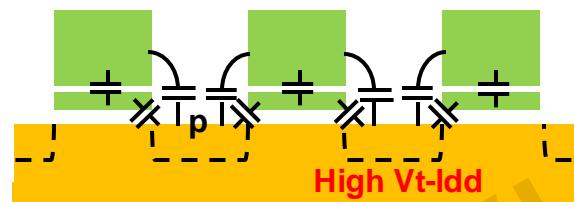
No LDD - DHC

1. With n- LDD (2D)



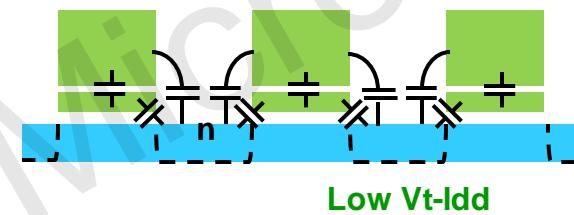
- LDD region has V_t which shifts by Q_{fg} and V_{wl_adj} .
- $V_{t_cell} = \max. (V_{t_ch.}, V_{t_ldd})$.
- We want $V_{t_cell} = V_{t_ch}$. Instead of V_{t_ldd} . V_{t_ldd} is more influenced by the neighboring cell and WL.

2. No LDD, P-type channel



- Therefore, lower V_{t_ldd} is a requirement, leading to n-type doped thin channel.
 - N-type is to reduce V_{t_ldd} .
 - Thin channel is to eliminate bulk leakage path to realizing n-type channel.
 - Thin channel also makes the V_t less sensitive to the trap density in the channel (less volume \rightarrow less charge)
 - Due to the poly Si trap (acceptor type trap), a net channel doping is ~un-doped with $1\sim3E18$ n-type poly.

3. No LDD, N-type thin channel



Vpass_R Required

$$V_{\text{pass_R}} \text{ Required} = \text{Average L7 } V_t + V_{BL} + DCR/DIBL * (2/3 * V_{BL}) + \gamma * \{\sqrt{1+V_{BL}} - \sqrt{1+1/3 * V_{BL}}\} + (V_{gs} - V_t)$$

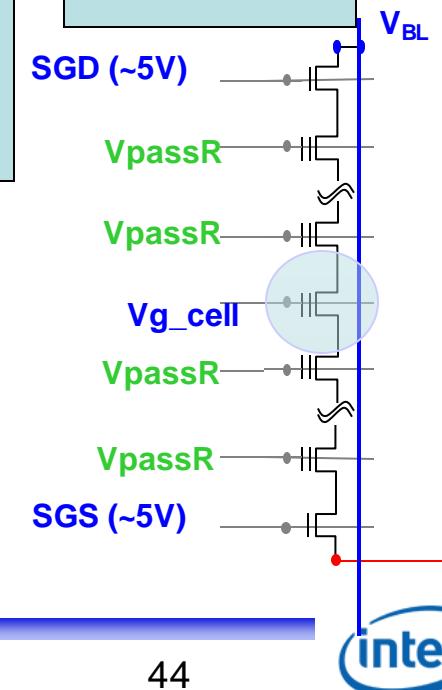
The V_t of the deselected cells. Including all sources, e.g. interference as well.

Since the sources of high WL#s is at $\sim V_{BL}$, V_g needs to be higher by this much

Going from PV7 verify to Pass gate condition, the V_{ds} across the high WL cells goes from $2/3 * V_{BL}$ to ~ 0 . So, the V_t increases

Going from PV7 verify to Pass gate condition, the V_s for the high WL cells goes from $\sim 0-1/3^{\text{rd}}$ VBL to VBL. So, the V_t increases by the body effect

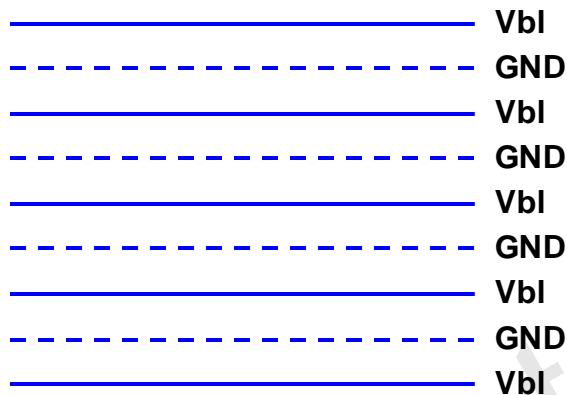
Gate overdrive to have acceptable string current



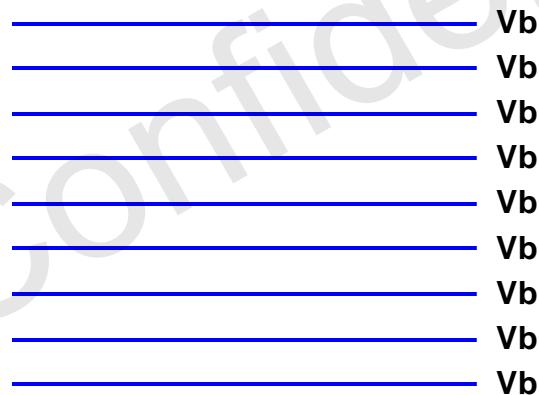
Higher $V_{\text{pass_R}}$ will help increase the string current, but will result in higher Read Disturb. So, there is Read Disturb to RWB trade-off to consider. (We will discuss Read Disturb later)

Shielded BL (SBL) and All BL (ABL)

SBL



ABL

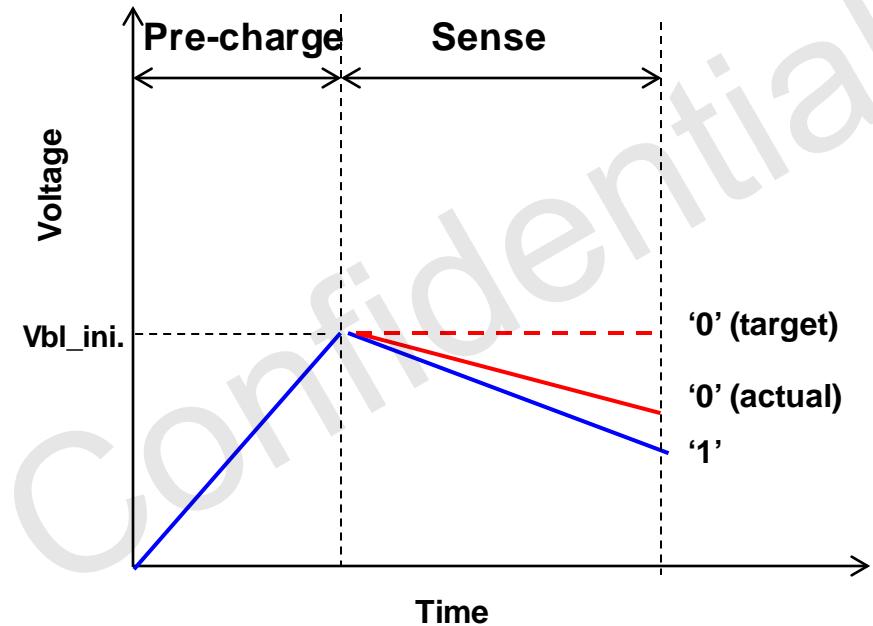
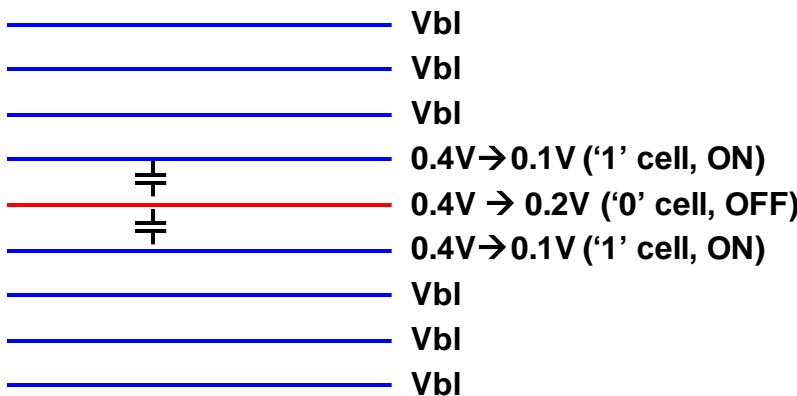


- Every other BL is selected for pgm/read.
- Page size is a half of # physical BLs

- All BLs are selected for pgm/read.
- Page size is the same as # physical BLs

ABL architecture is introduced to enhance pgm/read throughput by doubling a page size.

BL-BL Coupling Noise in ABL

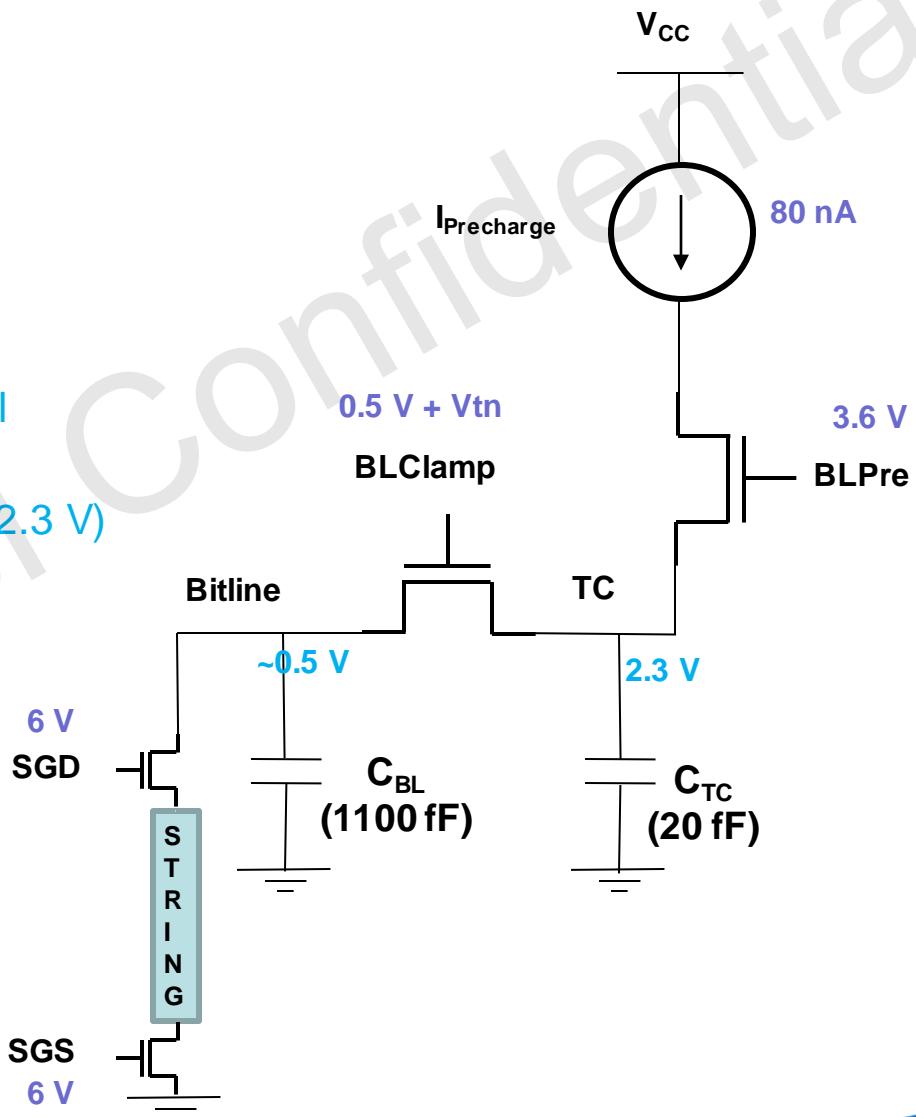


- '1' cells will draw current and discharge the BLs connected to them
- BLs connected to '0' cells are expected to stay at the pre-charge voltage
- Due to the BL-BL capacitance (~0.4pF/side), '0' BLs adjacent to '1' BLs are pulled down as well and discharges the sensing node.
 - '0' Cells are read as '1' → Sensing error.
 - Voltage sensing can't be used in ABL.
 - BL voltage has to be held constant during sensing → Active sensing

Active-Sensing – Step 1

Step 1: BL pre-charge/stabilization

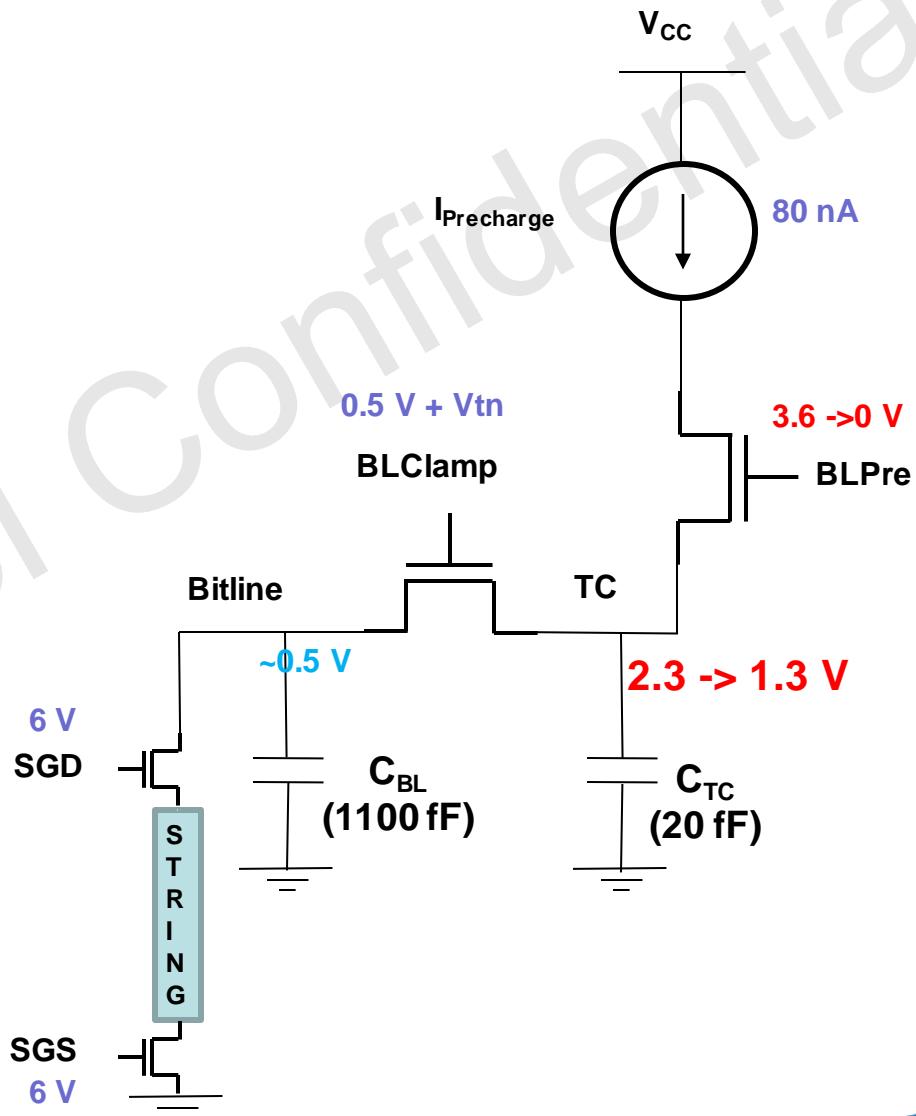
- Supply precharge current (80 nA)
- BLPre strongly ON (3.6 V)
- BLClamp targeting V_{bl} (0.5 V)
- SGS/SGD ON (~6 V)
- For strings drawing < 80 nA,
 - BL charges to approx 0.5 V. Exact V_{bl} depends on $I_{BLClamp} = I_{string}$ condition.
 - $I_{precharge} - I_{BLClamp}$ charges TC to V_{cc} (2.3 V)
- For strings drawing > 80 nA,
 - V_{bl} will settle at < 0.5 V with exact V_{bl} depending on $I_{BLClamp} = I_{string} = I_{precharge}$
 - Since all the precharge current is being discharged by string, TC cannot charge up to V_{cc} and is $\sim V_{bl}$
 - TC is lower than sense-amp flip voltage (~1.5V), and these will be sensed as '1' regardless of rest of sensing sequence.



Active-Sensing – Step 2

Step 2: Turn OFF BLPre

- BLPre OFF(0 V)
- Rest same.
- Since BLPre is OFF, there is no current feeding TC.
- TC will start getting discharged by ISTRING (= IBLClamp)
- If TC discharges below sense-amp flip voltage in a certain time, it is sensed as a '1' else it is sensed as a '0'.
- TC Develop time
 $= C_{TC} * (V_{CC} - V_{SA_flip}) / I_{sense}$
 $= 20 \text{ fF} * 1 \text{ V} / 20 \text{ nA} \sim 1 \text{ us}$
- As long as $TC > V_{BL} + \sim 0.2\text{V}$, V_{BL} doesn't change through develop time.



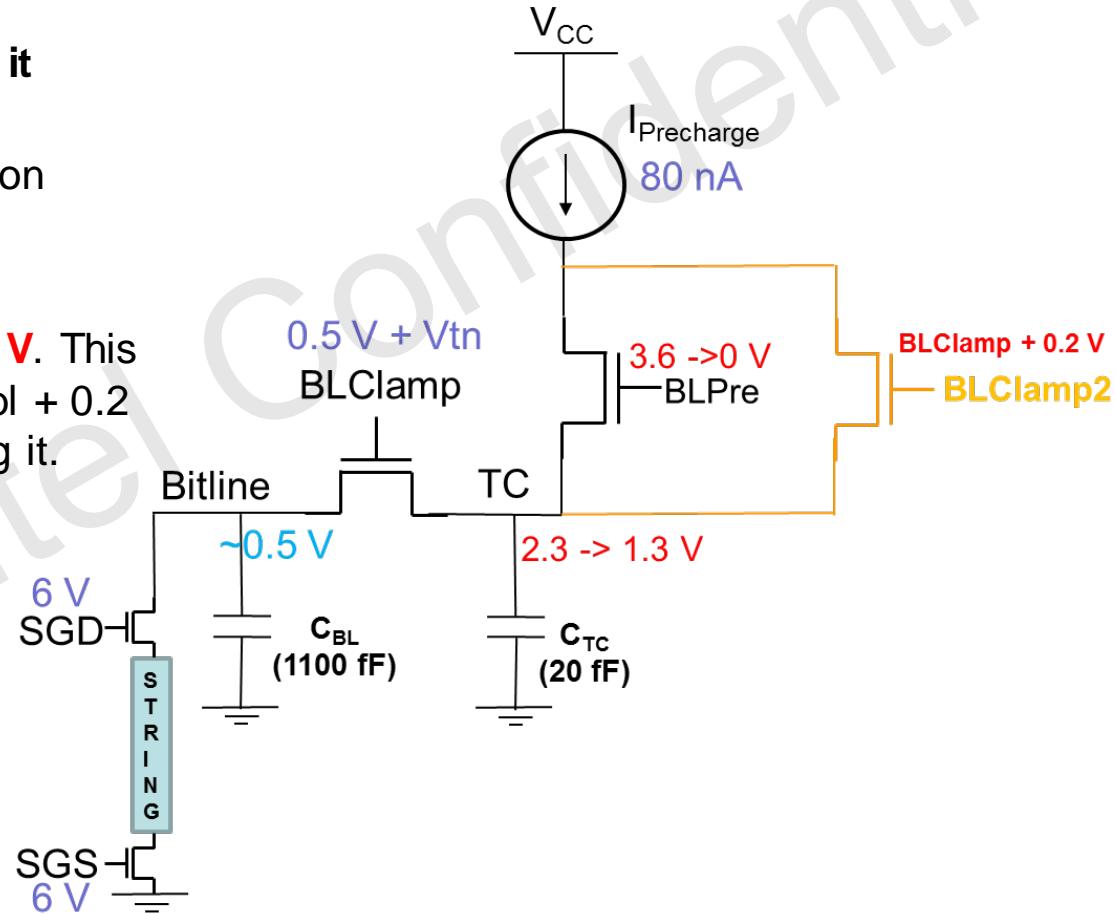
Active Sensing – Step 2

- How to prevent Vbl movement if TC discharges too much?

- Solution: Feed current into TC if it discharges too low.

- Using BLClamp2 (new device on active-sensing)

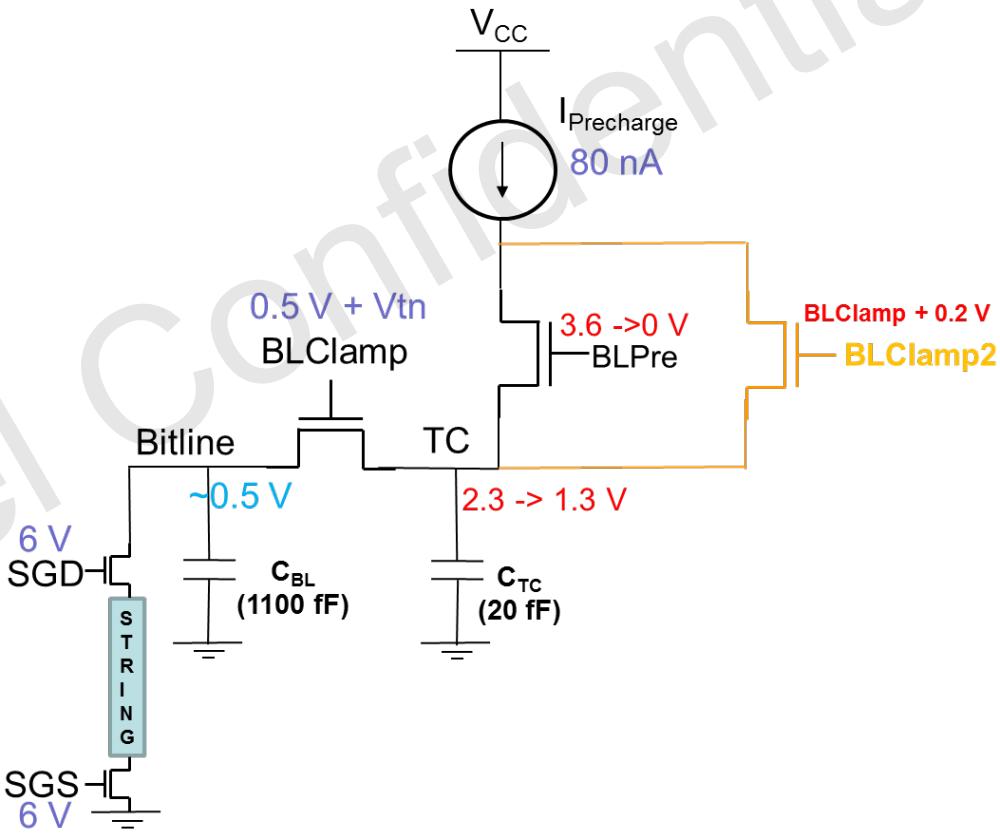
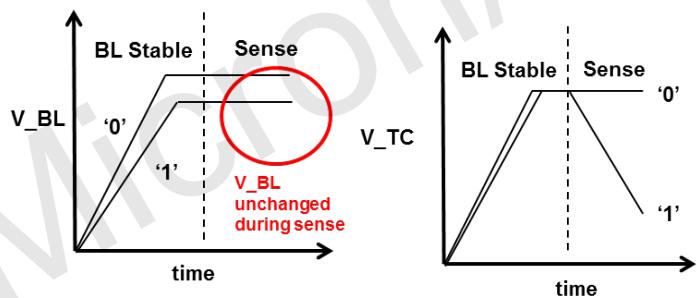
- Trim BLClamp2 to BLClamp + 0.2 V. This way TC never discharges below Vbl + 0.2 V because $I_{\text{precharge}}$ will start feeding it.



Active Sensing – Summary

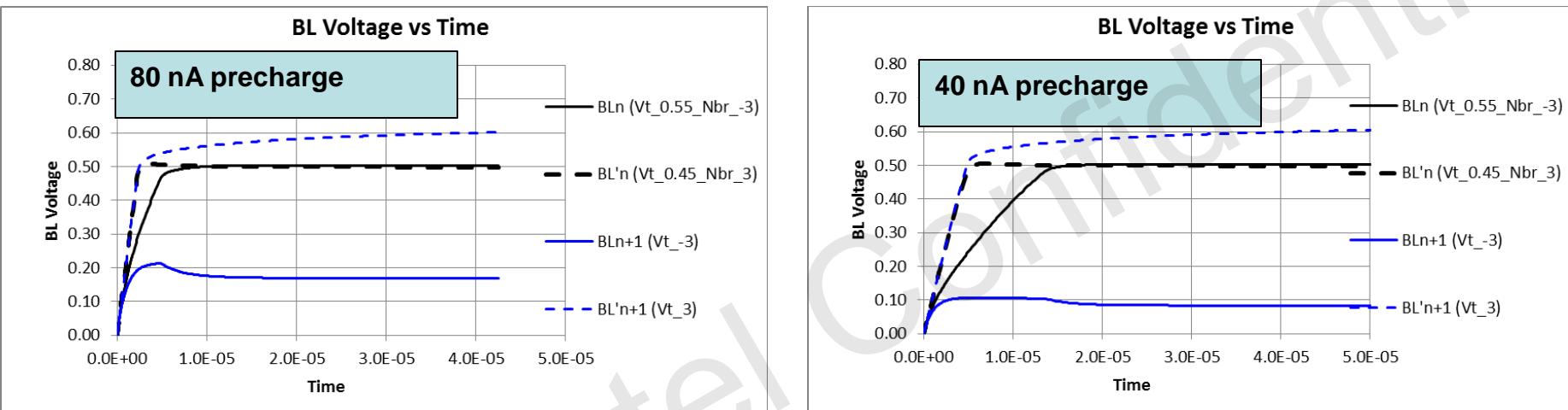
To summarize:

- Step1: Stabilize BL voltages (with SG ON) and precharge TC to Vcc.
- Step2: Turn off precharge current, and discharge TC. TC voltage is sensed after develop time
 - BL voltage doesn't change during TC develop time. This is accomplished by putting safeguard current into TC if it discharges too low, using BLClamp2.



Need for knockout (KO)

- Step 1 (BL stabilization) time depends on Iprecharge
 - Larger the Iprecharge, faster the BL stabilization.

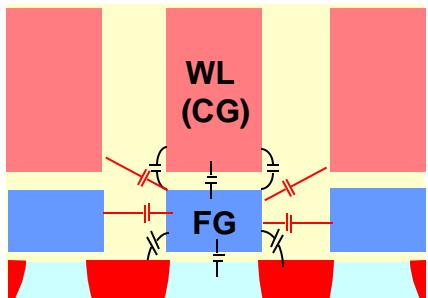


- Low Vt strings will be drawing full Iprecharge through entire sensing sequence.
 - Worse Icc / power
 - ‘SRC bounce’ leading to incorrect sensing.
- What is Source bounce?
 - Low-Vt cells draw large current constantly → IR drop in the SRC plate/line → Local SRC goes higher than global SRC (depending on location (R) and # low-Vt strings (I)).
 - If SRC is higher than intended, effective Vt looks higher.

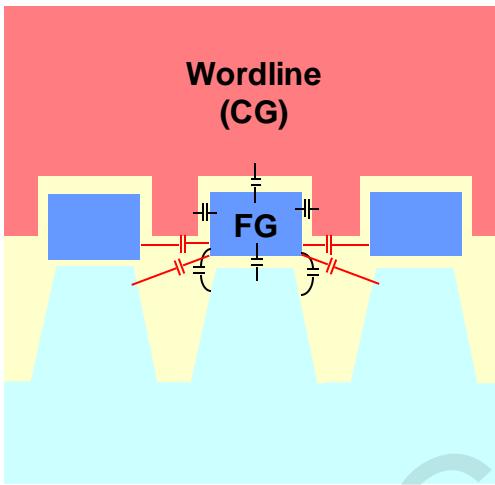
Knockout (KO)

- Knockout involves repeating the sense a second time, after knocking out the high current strings.
 - Coarse sense: Do full regular sensing sequence. Inaccuracy due to source bounce will make some Vt look higher. Use this sense to identify really low-Vt strings. Develop time decides target current.
 - Knock the low-Vt strings out (i.e. make their precharge current zero), then wait for BL stabilization (with BLpre/BLC clamp2 high) for neighbor BL/TC to recover.
 - Do a fine sense (take BLpre/BLC clamp2 low and develop TC again)
- Source bounce risk is higher for active-sensing because low-Vt strings are constantly drawing full precharge current whereas
 - Voltage sensing: No precharge current during develop.
 - Current sensing: Equal to sense current target ~ 20 nA, which is much smaller.
- Use of KO depends on tradeoff between
 - Using high Iprecharge to speed up BL stabilization, and taking penalty of KO/extrasense time.
 - Using lower Iprecharge to avoid KO/extrasense time but longer BL stabilization.

Cell to Cell Interference



2D

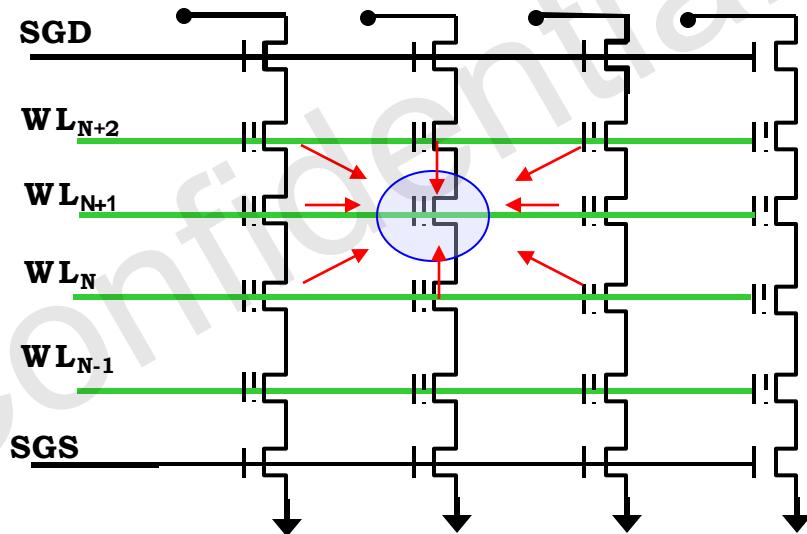


3D

- The Cell V_t is primarily decided by the charge on its floating gate.
- But the charge and voltage on the neighboring nodes (neighboring floating gate, Wordlines, channel, etc.) also impact the cell V_t.
- V_t of the victim cell can change due to the influence of the aggressor node onto the victim cell's Channel/LDD or the floating gate.
- Two key terms are:
 1. “FG-FG”: $\delta V_t_{victim} / \delta V_t_{aggressor}$
Aggressor can be BL, WL, or Diagonal neighbor → WL-WL, BL-BL, Diag
 2. “WL-FG”: $\delta V_t_{victim} / \delta V_g_{WL_{n+/-1}}$

FG-FG & WL_{n±1}-FG Coupling

- Impact of FG-FG interference:
 - During Programming all cells first start off in the Erased State.
 - Some cells are placed during a programming operation. Other cells (neighbors) are placed during a subsequent programming operation.
 - Effective V_t of the Cells Programmed first, moves higher as a result of the programming of the neighboring cells due to coupling from those cells.
- Impact of WL-FG Coupling:
 - Coupling from the neighboring wordline make the V_t depend on the neighbor WL voltage.



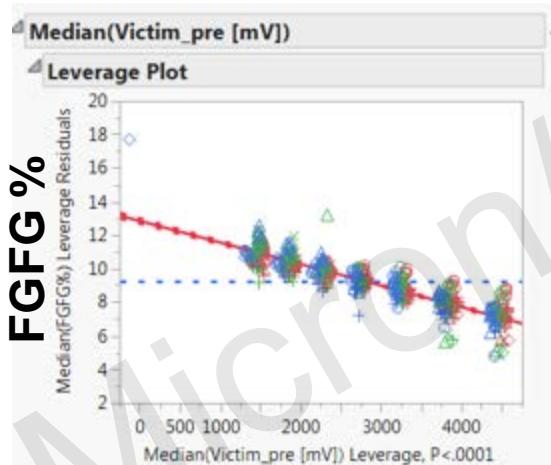
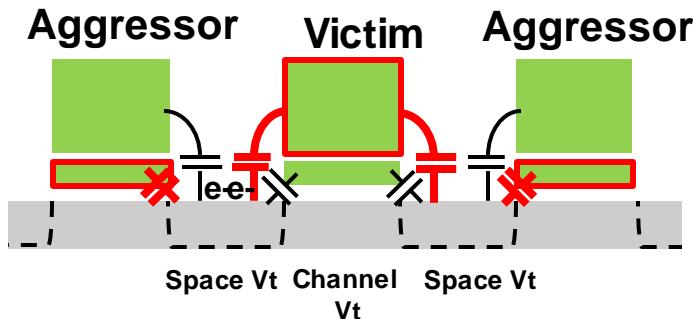
	80s	100s
FG-FG (WL-WL)	~4.5%	~5%
FG-FG (BL-BL)	~10%	~0%
FG-FG (Diag)	~1%	~0%
WL-FG	~10%	~13%

FGFG dependency on Vt level and cycling

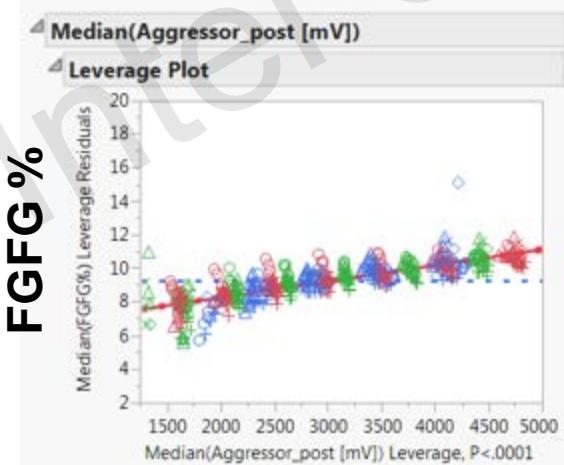
“FG-FG interference” is determined not only by FG-FG capacitance.

WL-Ch and FG-Ch play a role. too.
Space Vt has strong impact on ‘FG-FG’ Vt shift.

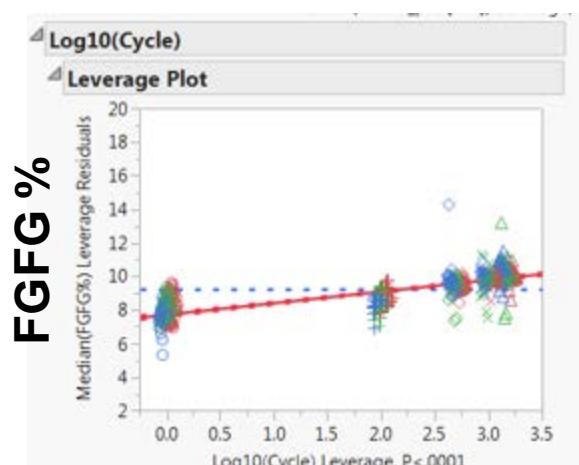
- Low Victim Vt (low WL voltage),
- High Aggressor Vt (more neg. FG potential),
- Cycling (trapping over space)
- Low temperature degrades ‘FGFG’ by making space Vt higher.



Victim Vt



Aggressor Vt



P/E cycling

MLC

MLC: 2 bits/cell: 4 states:

Level 0 = Erased

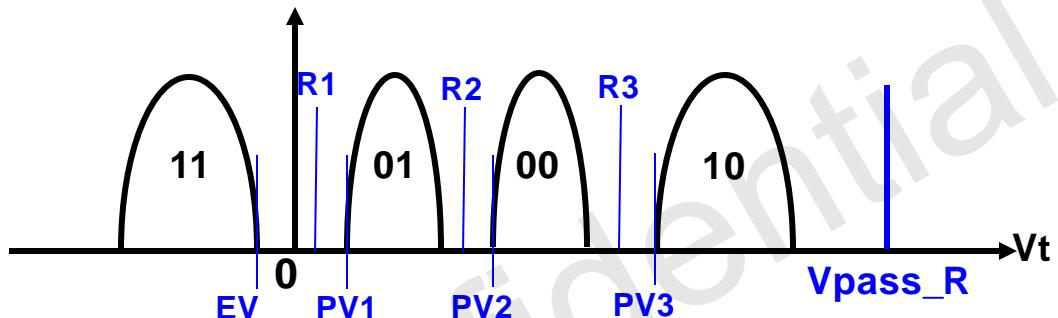
Level 1 = Pgm to V_{t1}

Level 2 = Pgm to V_{t2}

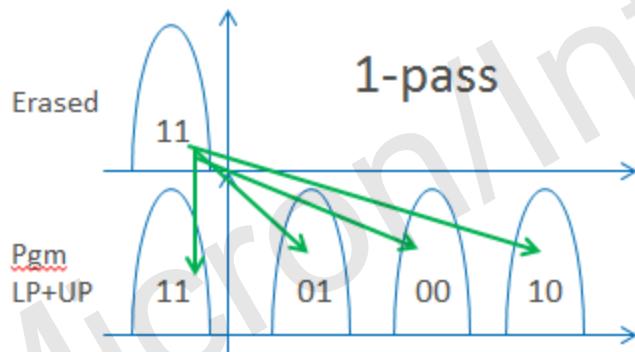
Level 3 = Pgm to V_{t3}

LSB (LP) Read with $V_g = R_2$

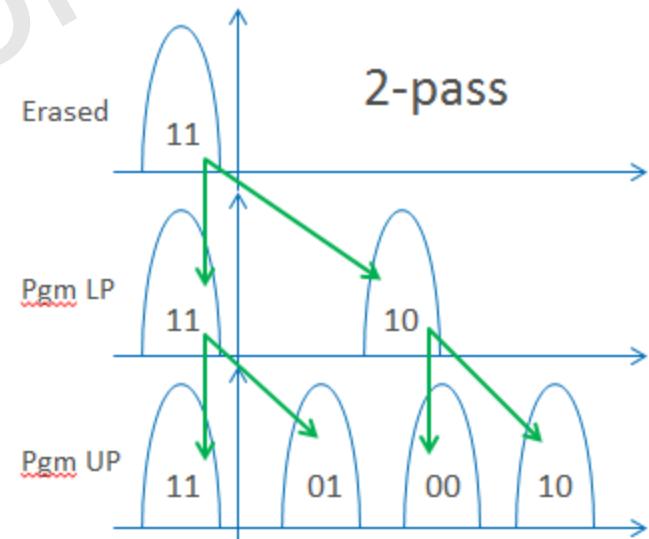
MSB (UP) Read with $V_g = R_1 \& R_3$



1 Pass vs. 2 Pass Programming



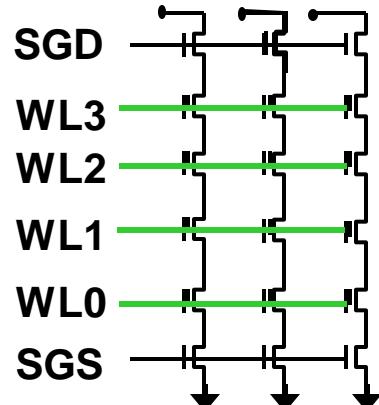
More susceptible to Cell-Cell interference.
More efficient for Program time



More immune to Cell-Cell interference.
Slightly inefficient for Program time

Back and Forth programming

Programming Order

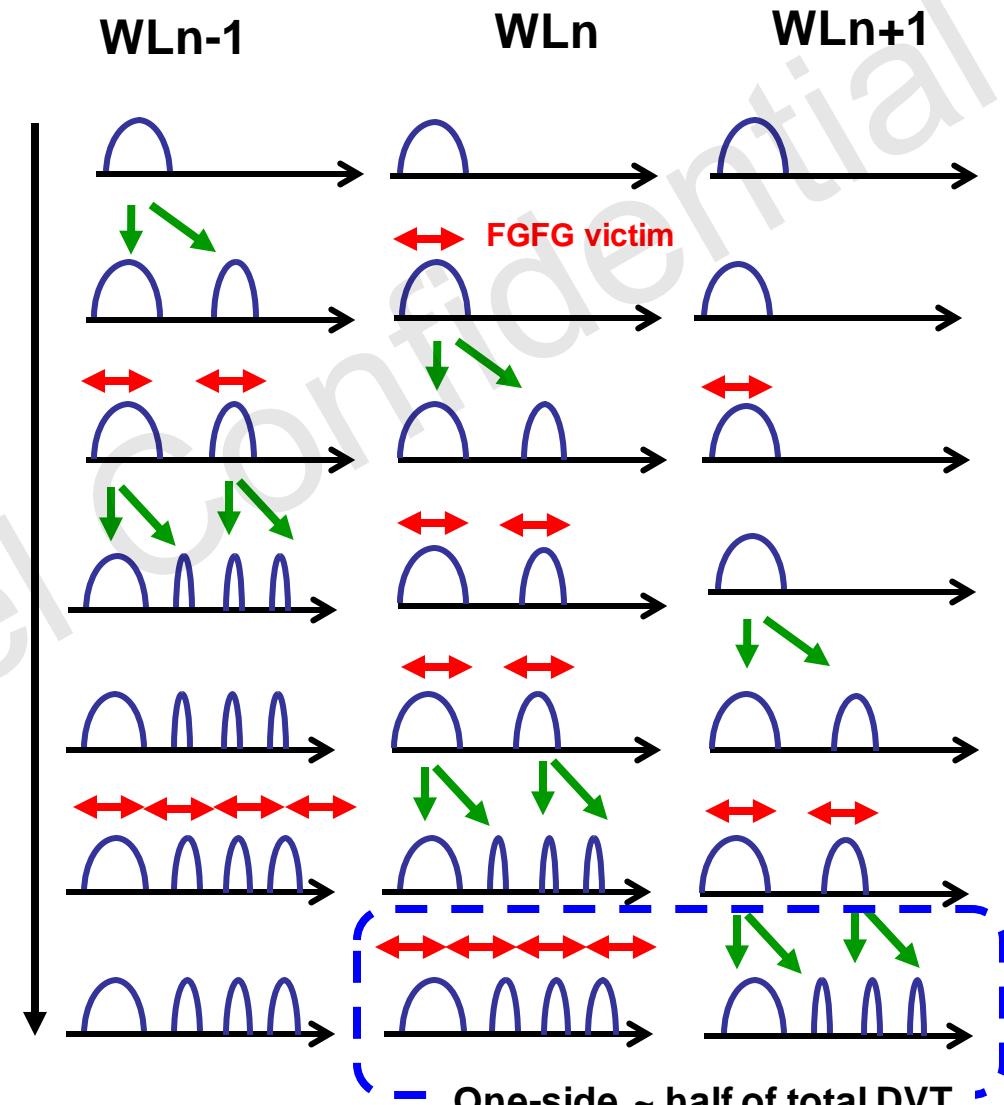


1st Pass (LP)	2nd Pass (UP)
5	7
3	6
1	4
0	2

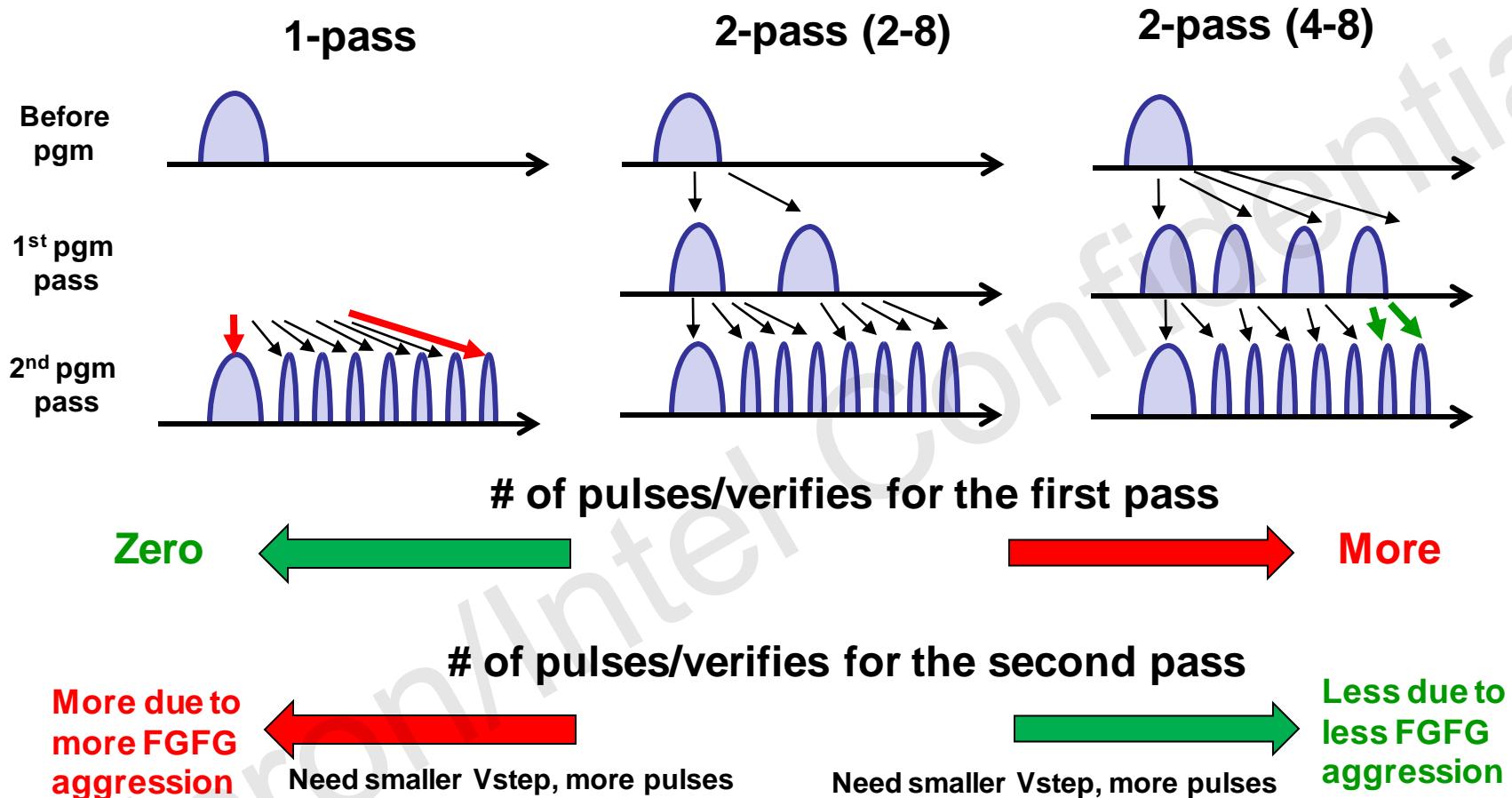
Programming is performed
From the source side to the drain side.

With multi pass programming,
Back and forth programming order.

This minimizes FGFG impact on the
Final Vt states. (1-side of WL, ~half of total DVT.)



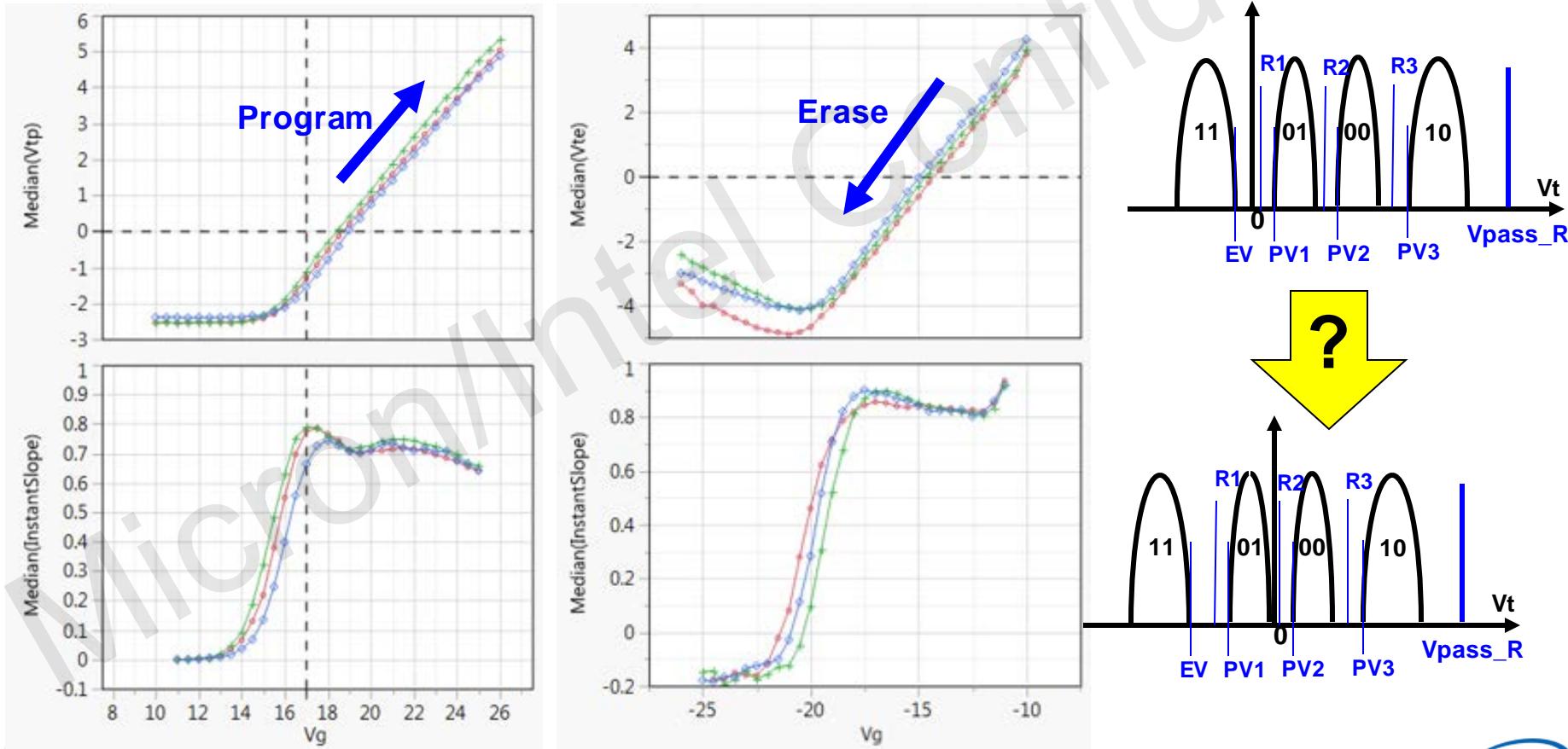
TLC programming sequence



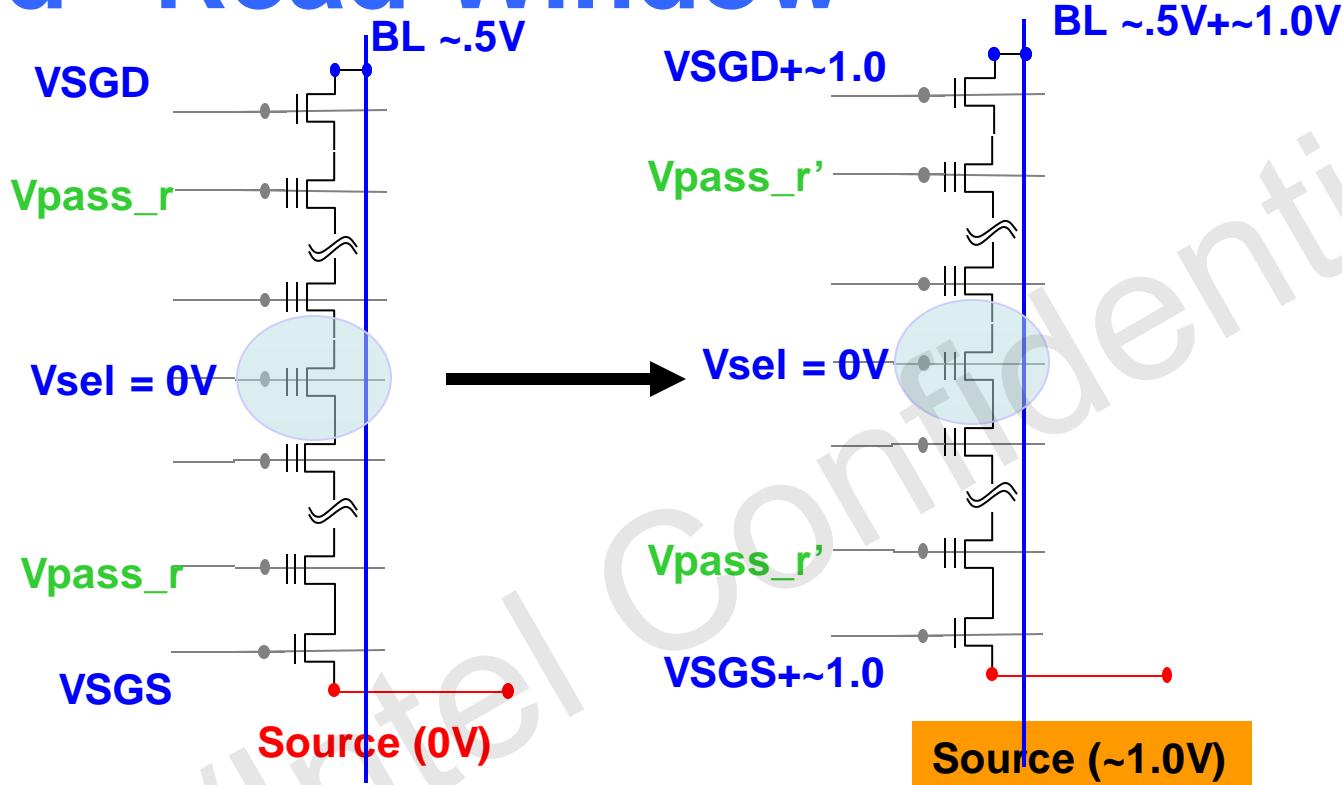
There are competing effect on tProg among programming algos.
The best algo. for tProg is FGFG dependent.
(FGFG<2%) 1-pass, (FGFG<4%) 2-8, (FGFG ~6%) 4-8.

“Shifted” Read Window

- Since NAND has somewhat symmetric Program/Erase window (e.g. +5V/-5V), ideally we would like to set the V_t window to be straddling the 0V.
- However, this requires ability to read -ve V_t → -ve WL voltage is required
- ➔ Adds design complexity and die size



“Shifted” Read Window

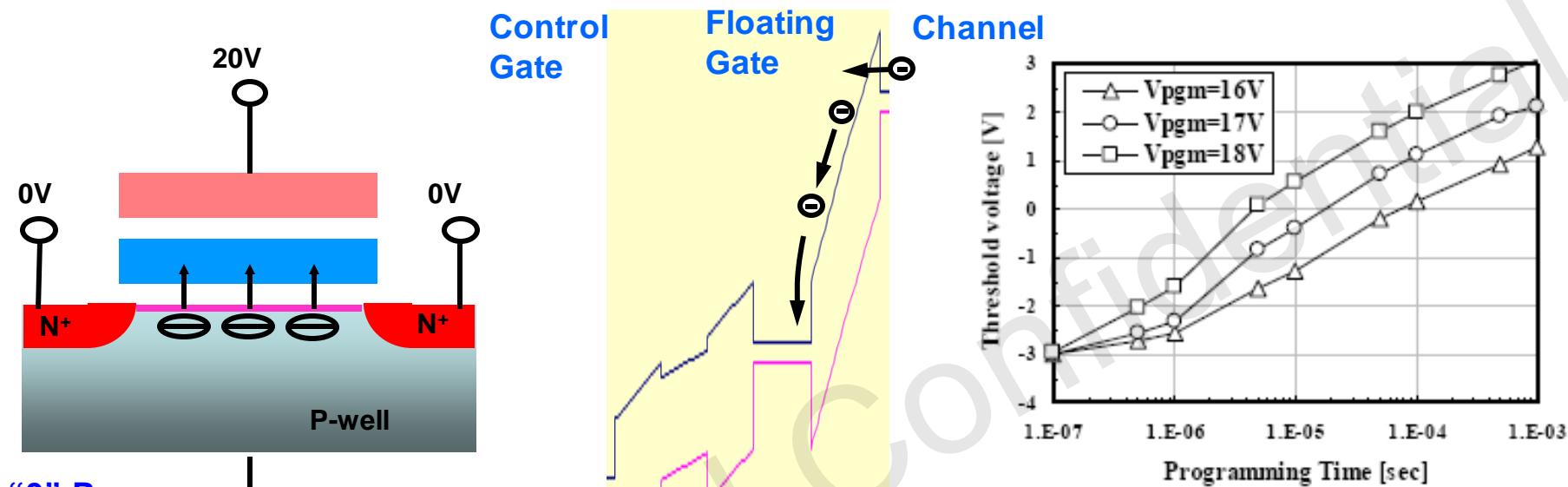


- We can enable –ve V_t read capability by shifting the source voltage +ve
 - With $V_s = 0V$ and $V_g = 0V$, select cell will not conduct if $V_t > 0V$ ($V_{gs} - V_t < 0$)
 - With $V_s = 1V$ and $V_g = 0V$, select cell will not conduct if $V_t > -1V$ ($V_{gs} - V_t < 0$)
- +ve source bias allows –ve V_t read capability without needing –ve voltage on die. Ability to shift window –ve is limited to 1.0-1.5V due to other circuit limitations.
- Shifted window allows another flexibility to optimize MLC/TLC V_t centering

Programming

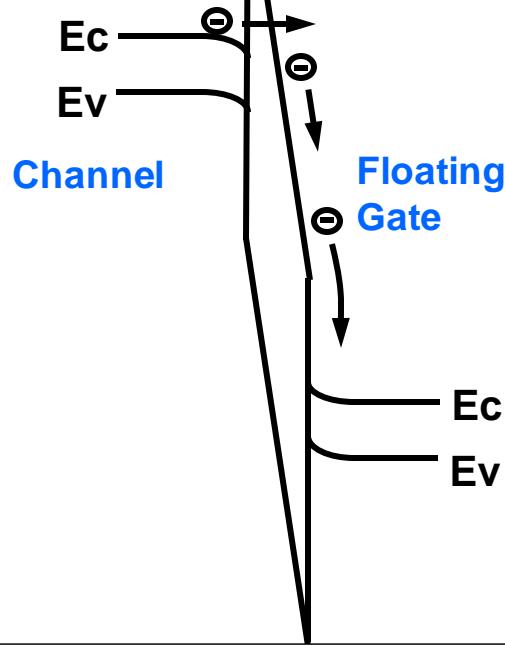
- FN Tunnel Programming
- Program/Erase Vt window requirement
- Program slope
- Program Saturation
- Curvature effect in 3D NAND
- Program Vt sigma
- Programming algorithm
- Program Noise
- SSPC / Proximity Effect
- Program Inhibit / Program Disturb
- Boost loss mechanisms
- Hot-e Disturb
- SG Leakage

NAND Flash Programming - FN Tunneling



- Tunnel Programming from channel by biasing the Top Gate positive with respect to the channel and S/D.
- Actual Program Time (cumulative pulse width) ~100us. Program current ~ Displacement and Tunneling current. Low current allows large parallelism

FN Tunneling

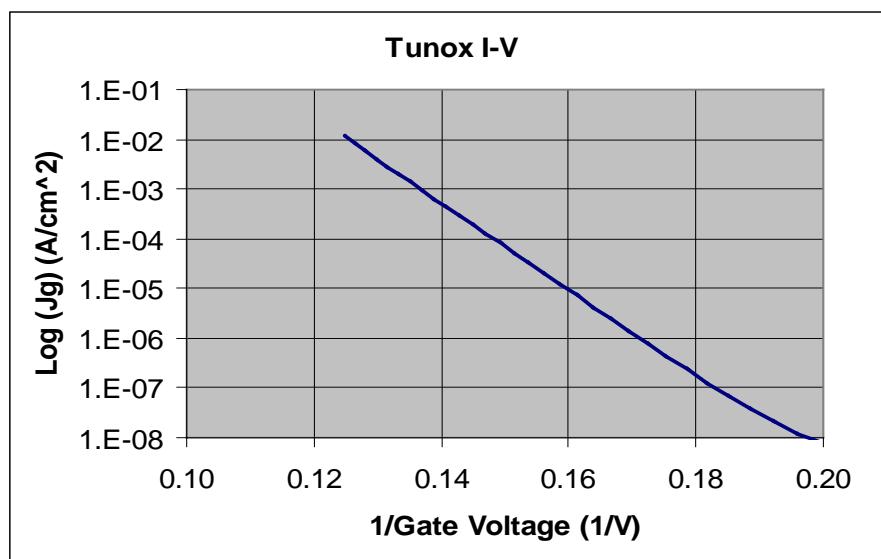
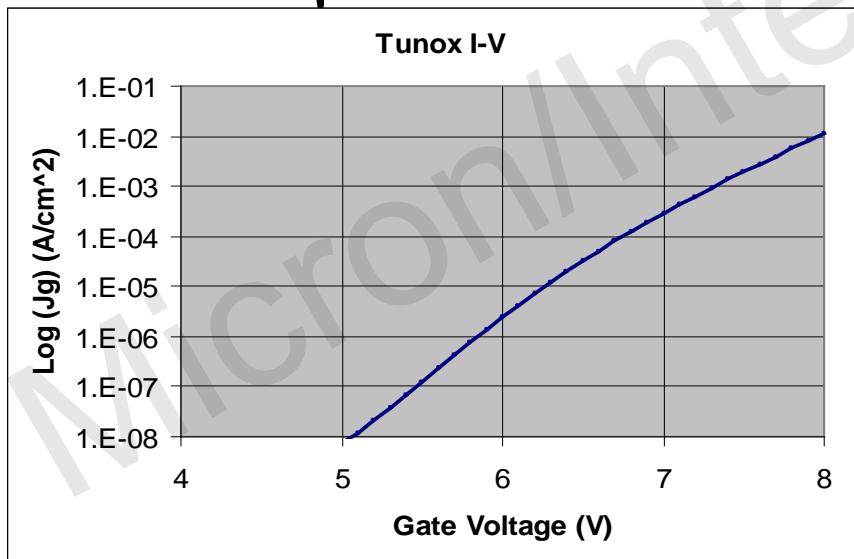


- Electrons tunnel through the first 25-30 Å of SiO_2 and then drift to the FG
- Fowler-Nordheim tunneling equation says that current varies with field as:

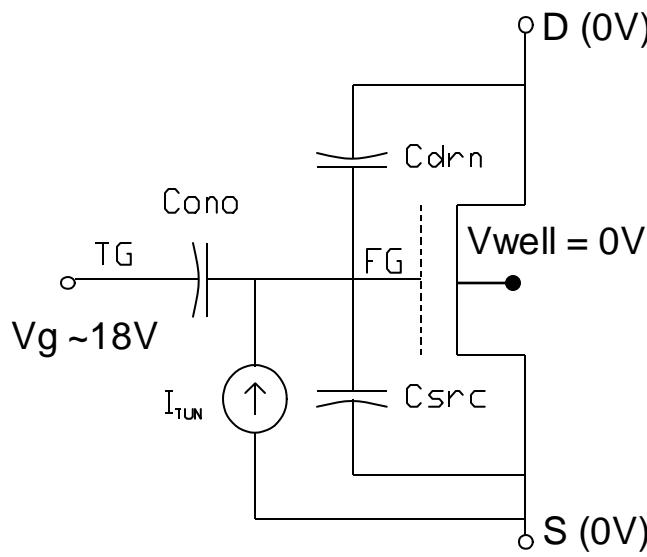
$$J = a E^2 e^{\left(\frac{-b}{E}\right)}$$

Coefficients a and b are also commonly known as α and β

- The current changes by ~10X with ~0.8V change in the FG voltage

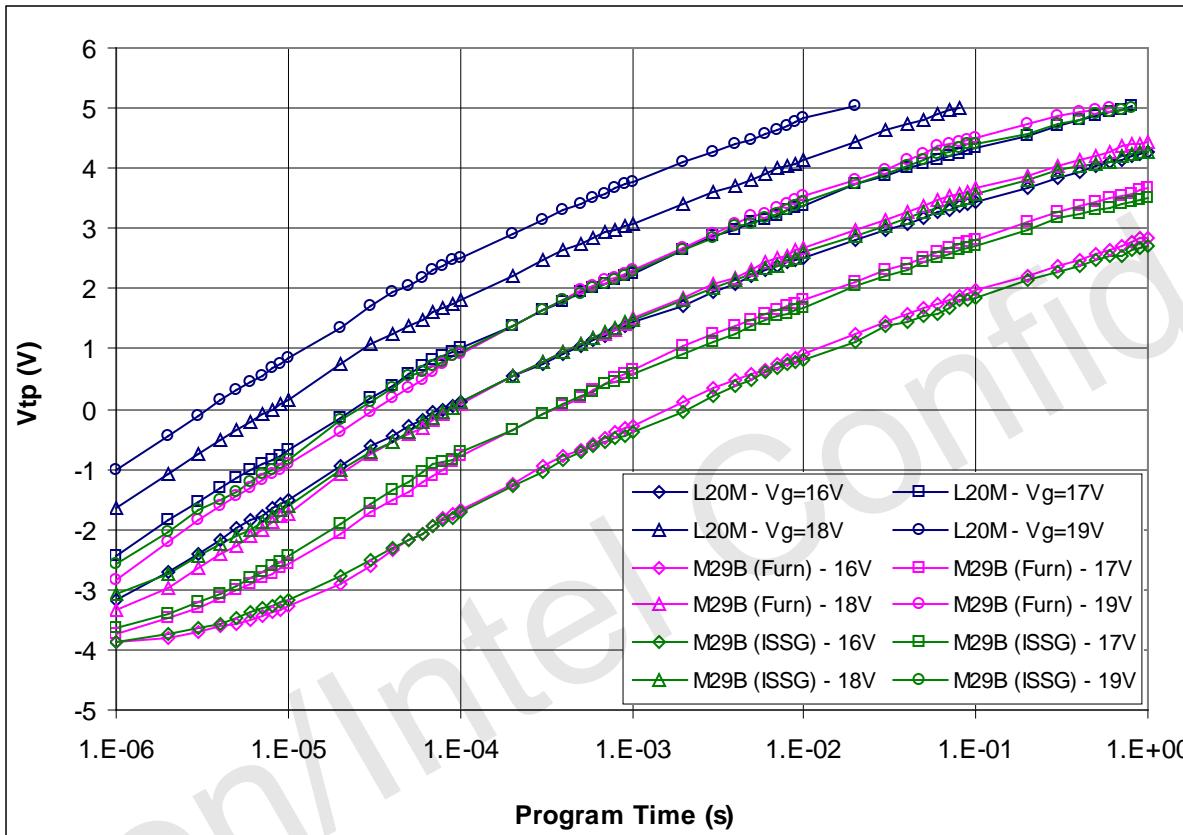


Program V_T vs. Time



- Equiv. ckt shows Floating Gate getting charged by I_{TUN} (actually –ve current as it is the electrons tunneling onto FG)
- At the beginning of program, V_{FG} is high and I_{TUN} is large, so high tunnel current and programming is fast
- As number of electrons on FG increases the FG voltage reduces (by $\Delta Q_{fg}/C_{tot}$) decreasing I_{TUN} as given by the FN Equation
- Each $\sim 0.8V$ decrease in V_{FG} drops I_{TUN} by $\sim 10x$
 - Programming rate drops by 10x, 100x, etc as the floating gate gets charged by .8V, 1.6V, ...
 - Makes Programming vary as a "log time"

Programming



Example Exhibit only.
Exact V_t vs V_g
dependence is very
technology specific
(Tunox thickness,
GCR, UV_Vt, ..)

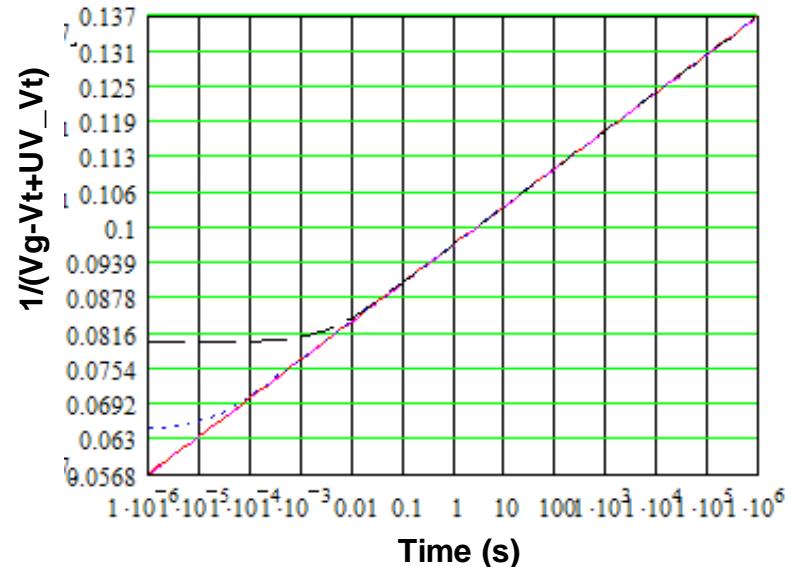
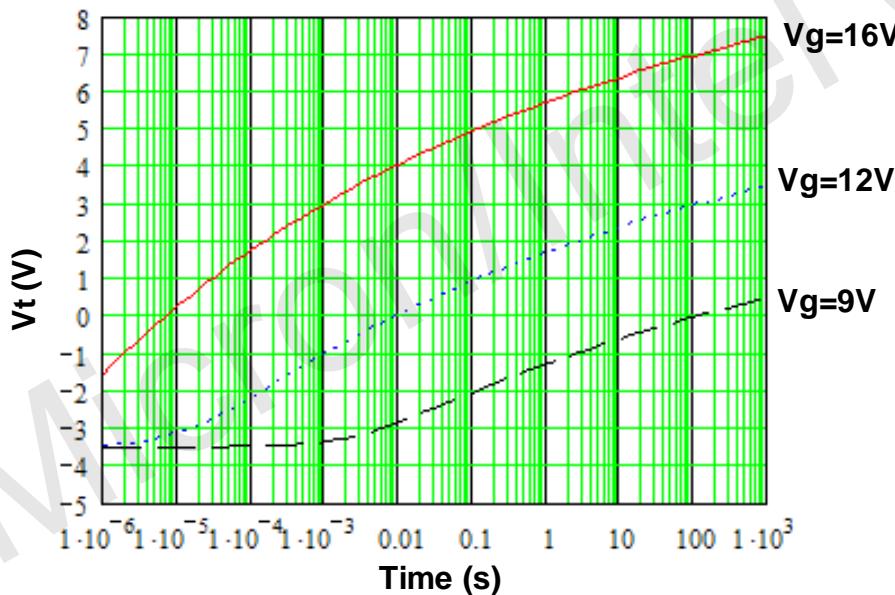
- Log Time behavior of the program V_t can be seen above
- A rough rule of thumb is that the cell V_t changes by 1.3-1.5V/dec under programming conditions (~10-100us). The slope decreases to ~0.5-0.6V/dec as we go longer in time (~30-300s).

1/(Vg-Vt) Curves

Plotting $1/(Vg - Vt)$ vs. \ln_t gives a straight line with an intercept and slope given by:

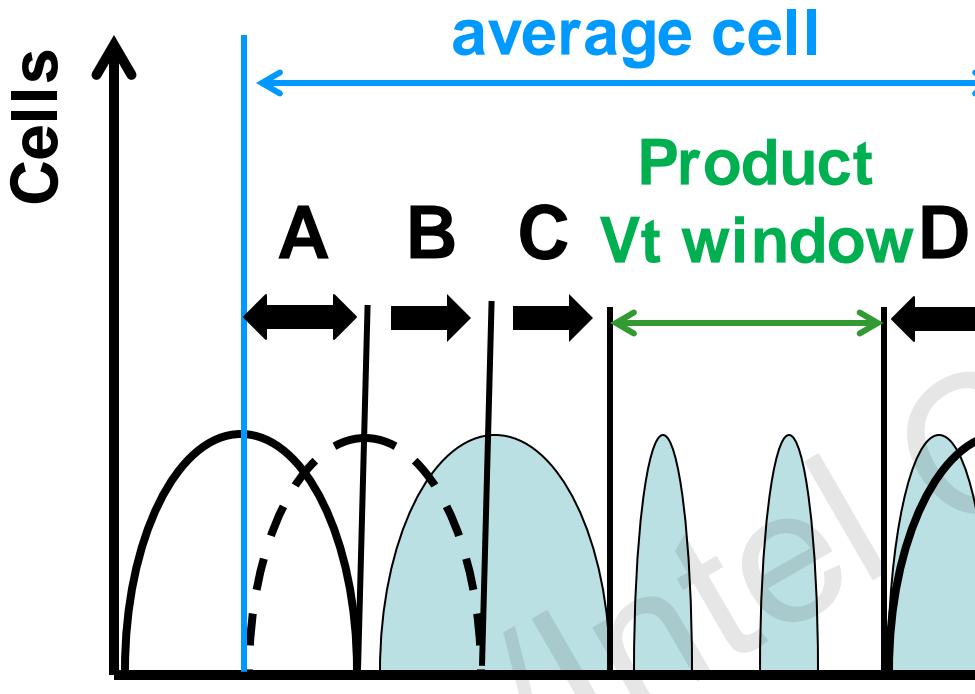
$$\left[\left(Vg - Vw + Vth \cdot Gcr^{-1} \right) - Vtx \right]^{-1} := \frac{Gcr}{\beta \cdot T_{ox}} \cdot \ln \left(\frac{\text{Area} \cdot \alpha \cdot \beta \cdot 1s}{C_{tot} \cdot T_{ox}} \right) + \frac{Gcr}{\beta \cdot T_{ox}} \cdot \ln \left(\frac{t}{1s} \right)$$

- Where Vg is the gate voltage, Vth/Gcr is the UV_Vt of the cell
- Vw is the Well or S/D voltage (depending on erase or programming)
- α and β are the FN tunneling coefficients. Areas is the tunneling area and C_{tot} is the Total FG capacitance. T_{ox} is the Tunnel-oxide thickness



P/E window requirement

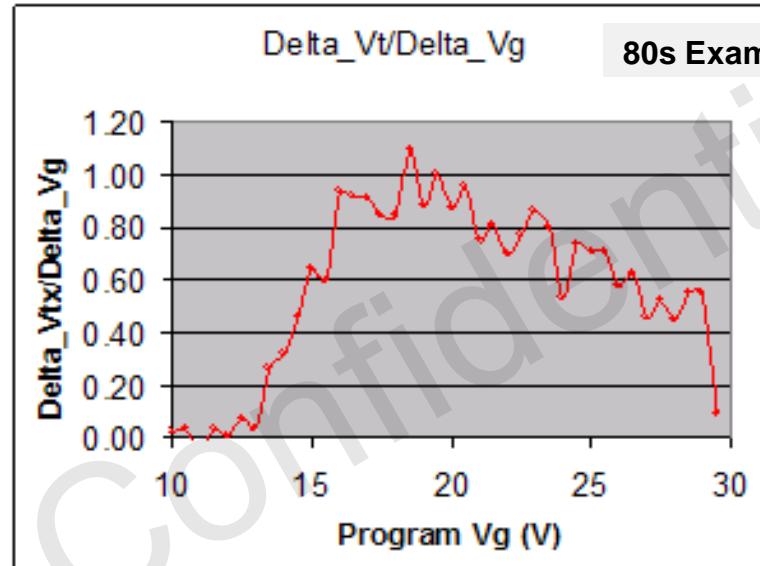
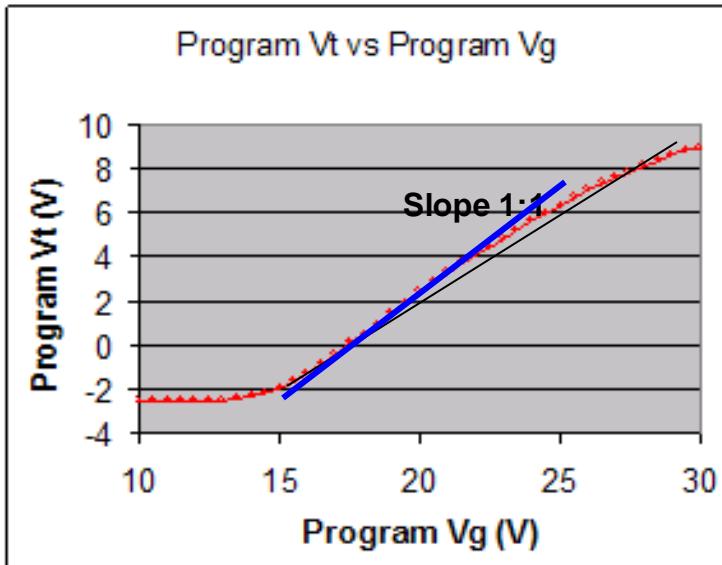
P/E Vt window of



- In order to support ~4V P/E window in product, Typical cell needs to have >8V P/E window.

Product P/E window	PV3-R1	~4V
A. Cycling trap up margin	V _{te} _sat. degradation post cycling	~2V
B. FGFG margin	Erased V _t shifts up after neighbor cells programmed	~0.25V (=10% x 2.5V (average V _t shift with RRP data))
C. EVS margin	Cell to cell V _{te} variation margin	~2V (=0.25V x 4sigma (~+1 sigma from ECC requirement x2 (EVS and PVS)))
D. PVS margin	Cell to cell V _{tp} variation margin	
Required Cell P/E window		> 8.25V

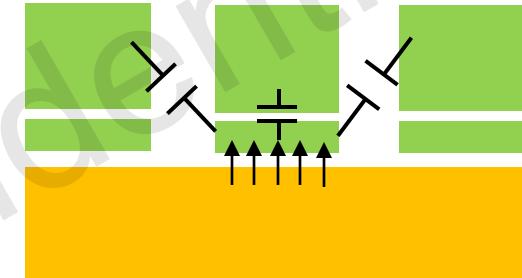
Pgm Vt vs Program Vg (80s)



- Programming is done as ISPP (incremental step pulse programming)
- For a given programming pulse width, the program Vt tracks the program gate voltage. $V_{gVt} (= V_{pgm} - V_{tx})$ is a good cell programming metric.
- V_{gVt} is primarily determined by the Cell UV_Vt, FN Tunneling coefficients, Tunnel-oxide and ONO thicknesses, GCR, etc.
 - $V_{gVt}-UV_Vt$ is proportional to Cell EOT ($= Tox/GCR$)
- We would expect $\Delta V_{tx}/\Delta V_g$ to be ~ 1.0 . However, some deviation from 1.0 is seen due to depletion effects (GCR difference between Program and Read), and CG impacting the Vt directly (instead of through FG).

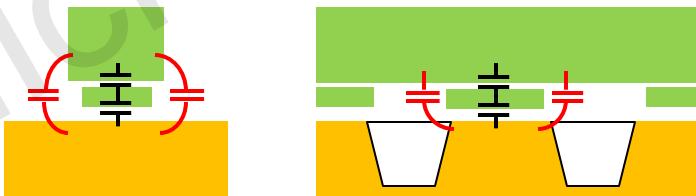
Programming Slope

- Typically the program slope should be = 1.0 (Assuming fixed Vinh).
- When you increase the V_{cg} (during Pgm) by ΔV_{cg} , the V_{fg} moves up by $GCR * \Delta V_{cg}$ and as a result we inject electrons into the FG and FG comes down by $GCR * \Delta V_{cg}$. Later when you read the cell V_t , you will need to apply a higher V_{cg} given by $GCR * \Delta V_{cg} / GCR = \Delta V_{cg}$ to get the cell to turn on.
 - $\Delta V_t = \Delta V_{cg}$ during programming! → Program slope = 1.0
- If the $Vinh$ on the adjacent WL was changed from pulse to pulse, then there will be coupling from the adjacent WL and which would have moved the V_{fg} by more than $GCR * \Delta V_{cg}$, so ΔV_t will be more
 - $\Delta V_t = \Delta V_{cg} + 2 * WL-FG * \Delta Vinh$
- If the GCR during programming and Read were different (e.g. due to Poly depletion effects and different IPD fields during Programming and Read):
 - Program slope = GCR_p/GCR_r (Assuming fixed $Vinh$)

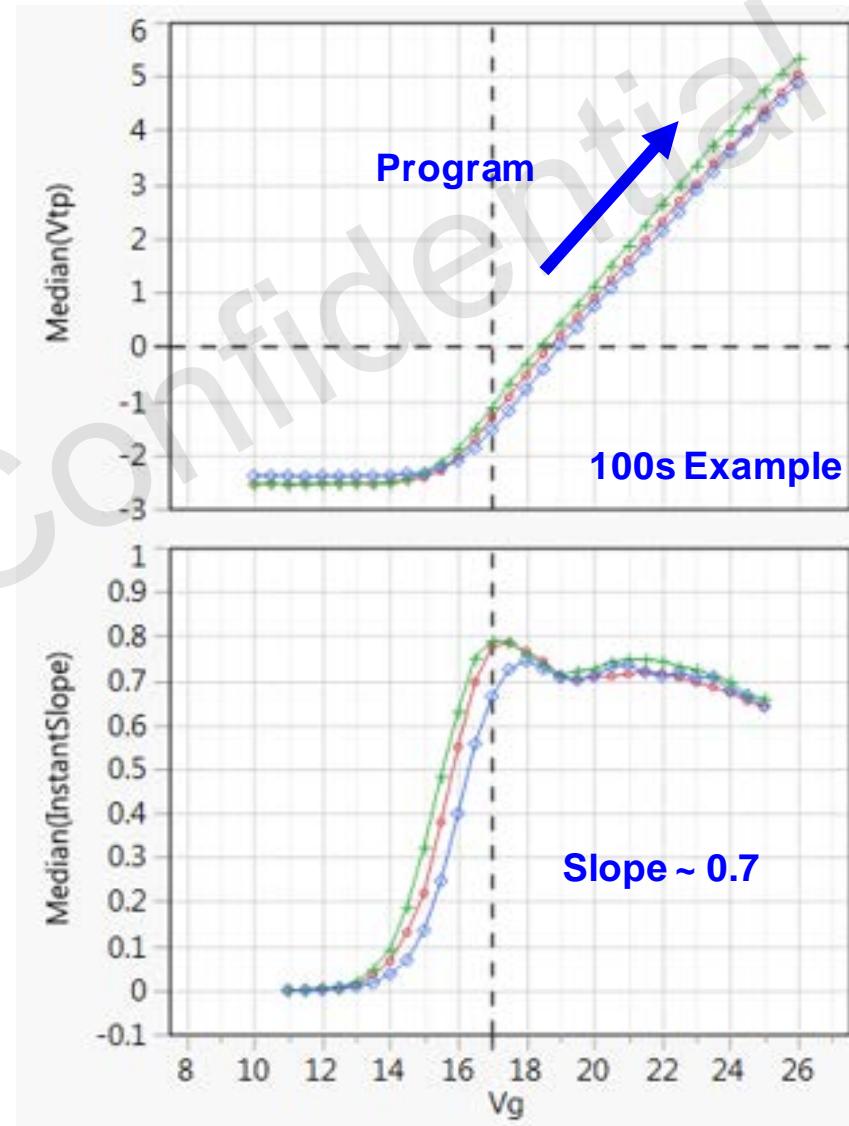


Programming Slope

- Additional factors impacting pgm. Slope:
 - As the cell programs, field across the IPD increases and there could be a competing effect causing increasing IPD leakage causing <1.0 Program slope.
 - In the limiting case, cell will stop programming (Programming saturation: V_{tp_sat}).
 - (Partial capture (e.g. 70%) of the electrons in the FG does NOT change the program slope. It can cause $V_g V_t$ to increase slightly, but slope will not be affected as long as capture % doesn't change)
 - Direct CG-AA capacitance can reduce Program slope:
 - Program slope $\sim \frac{C_{eo}}{C_{eo} + C_{CG-Ch}}$

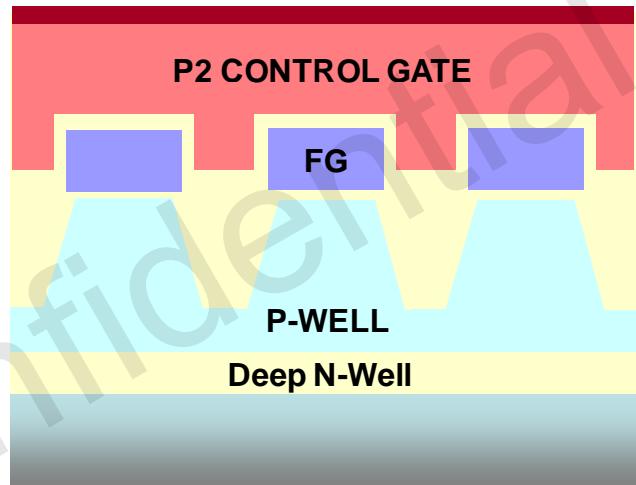


Direct CG-AA capacitance

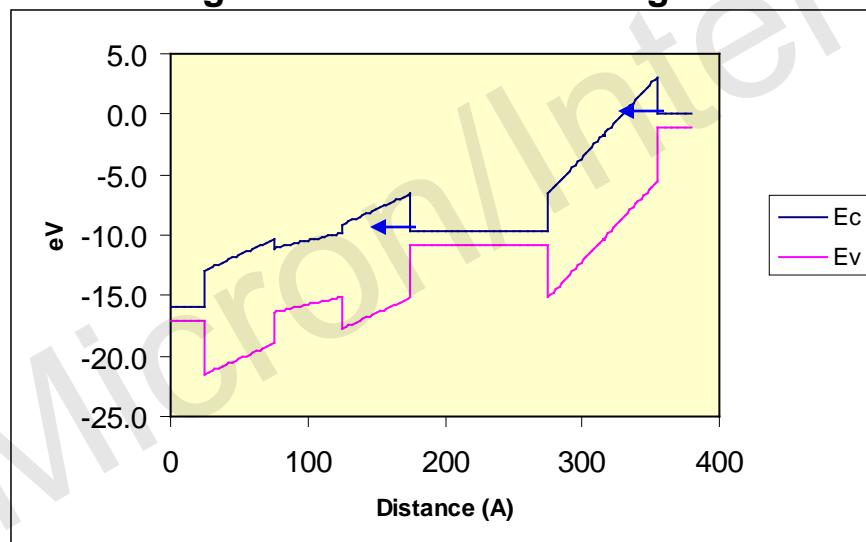


Program Vt Saturation

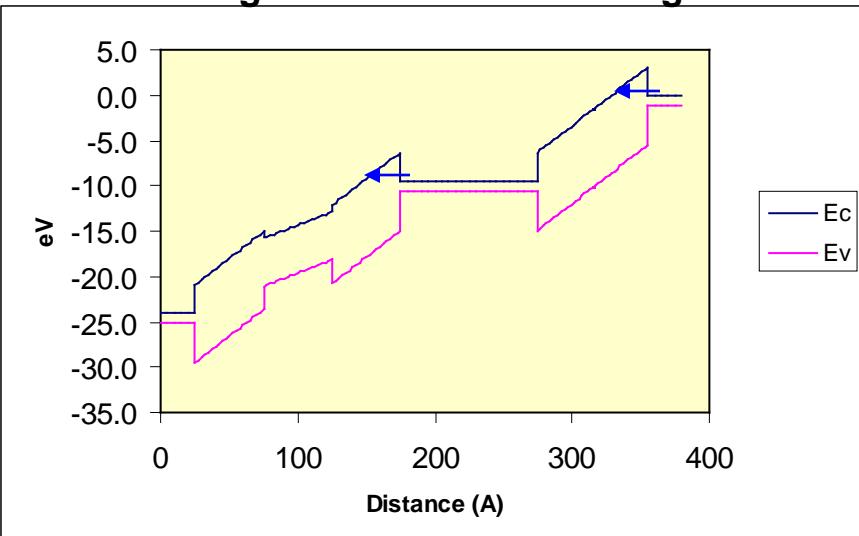
- Vt saturation at high Vt because of electrons tunneling out through the IPD.
- As the programming gate voltage increases and the cell gets programmed higher and higher, the field across the IGD increases and the leakage through the IGD starts becoming large and can become equal to the tunneling current



Low Program Vg and Vt ($V_g = 16V$, $V_t = 0V$).
Current through ONO << Current through Tunox

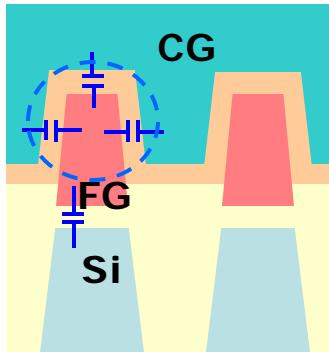


High Program Vg and Vt ($V_g = 24V$, $V_t = 8 V$).
Current through ONO = Current through Tunox

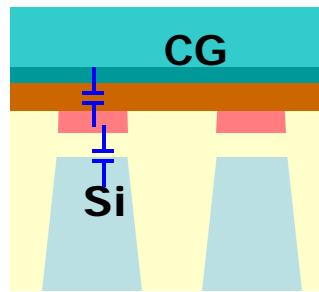


Tunox = 80A, ONO = 150A Physical / 125A Electrical, Wrap Ratio ~ 2.5

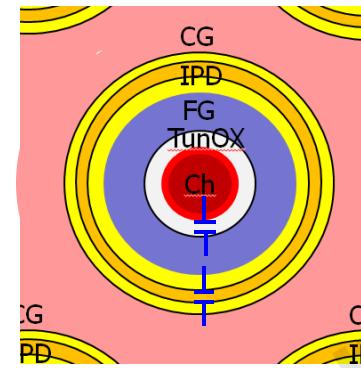
3D NAND Pgm/Ers Window



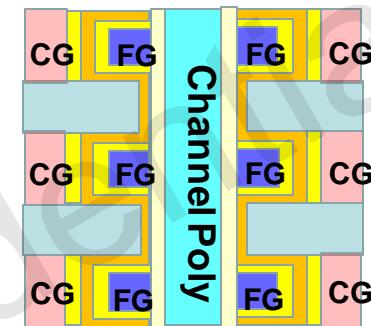
70s: Wrap



80s: High-K / MG



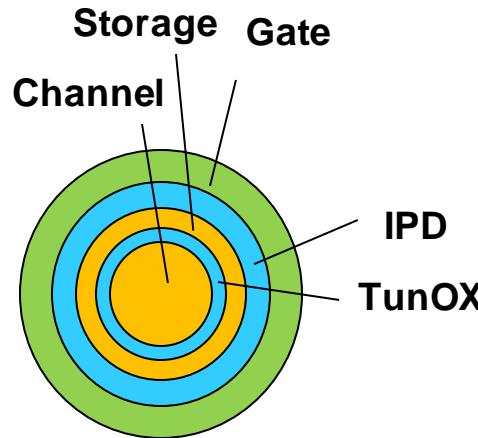
100s 3D cell



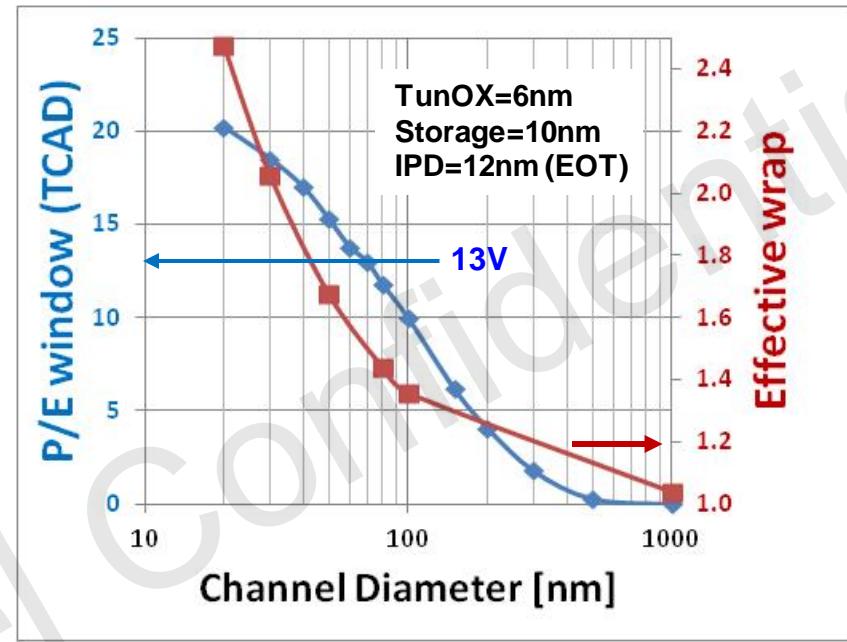
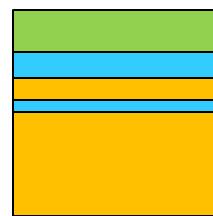
- Good Program/Erase window requires high Gate coupling Ratio (FG couples stronger to the Control Gate) → Higher Tunnel-oxide field and Lower IGD field
- 3D NAND (100s) cannot have conventional wrap (like 70s) or High-K/MG (like 80s). 3D NAND counts on the circular shape to provide higher coupling area between FG & CG relative to FG & Si.
- Alternatively it can be thought of as electrical field being higher closer to the center of the circle (e.g. Si/Tunox interface) than outside (FG/IGD interface) which allows FN current through Tunox but not IGD

GAA Wrap Effect

Gate-all-around (GAA)

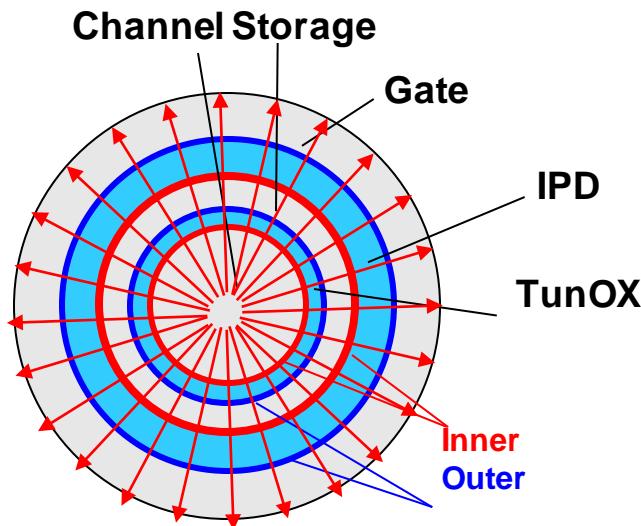


Planar

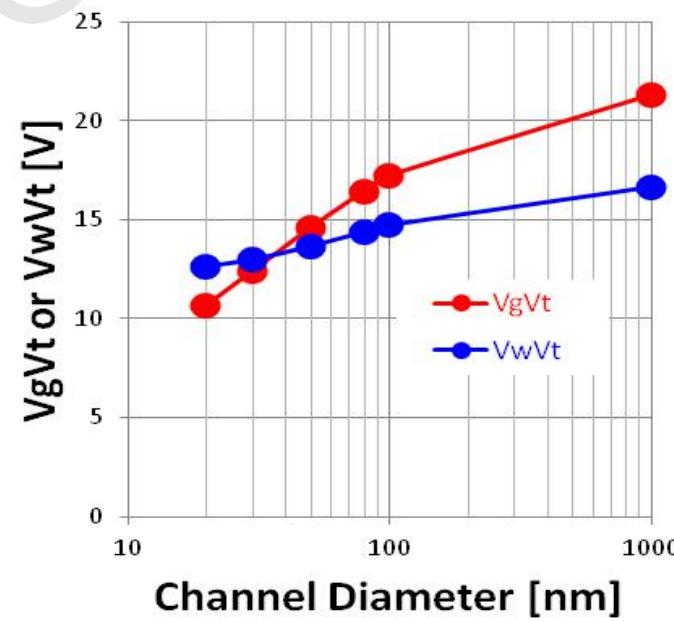
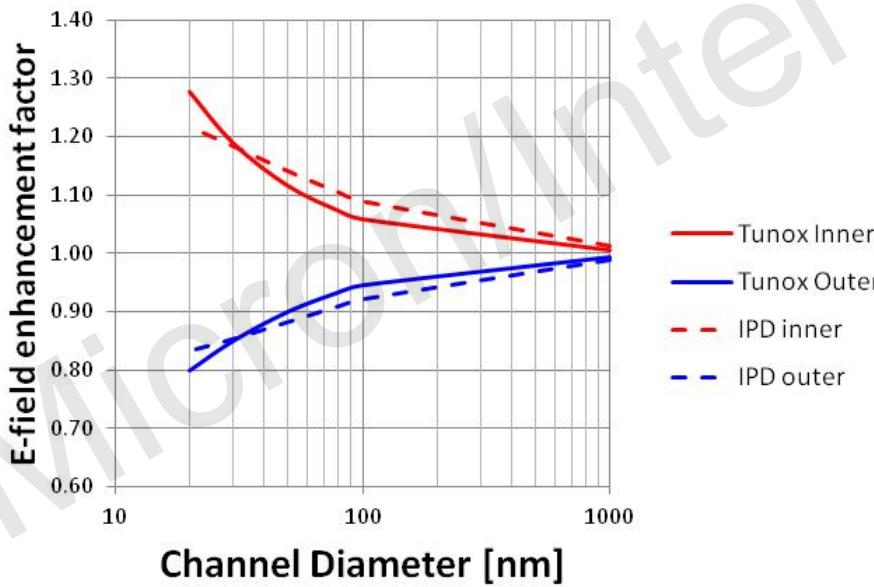


- Due to GAA architecture of 3D, IPD area is larger than TunOX area, providing effective wrap ratio > 1.0
- This increases ratio of $E_{field_tunox} / E_{field_ipd}$, enabling a wider program/erase V_t window.
- <70nm channel diameter is required to achieve >13V P/E window.

Curvature E-field effect

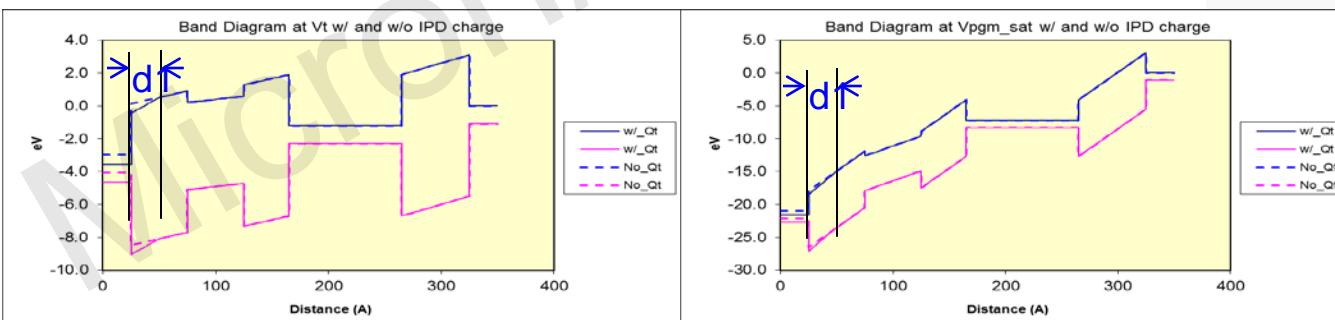
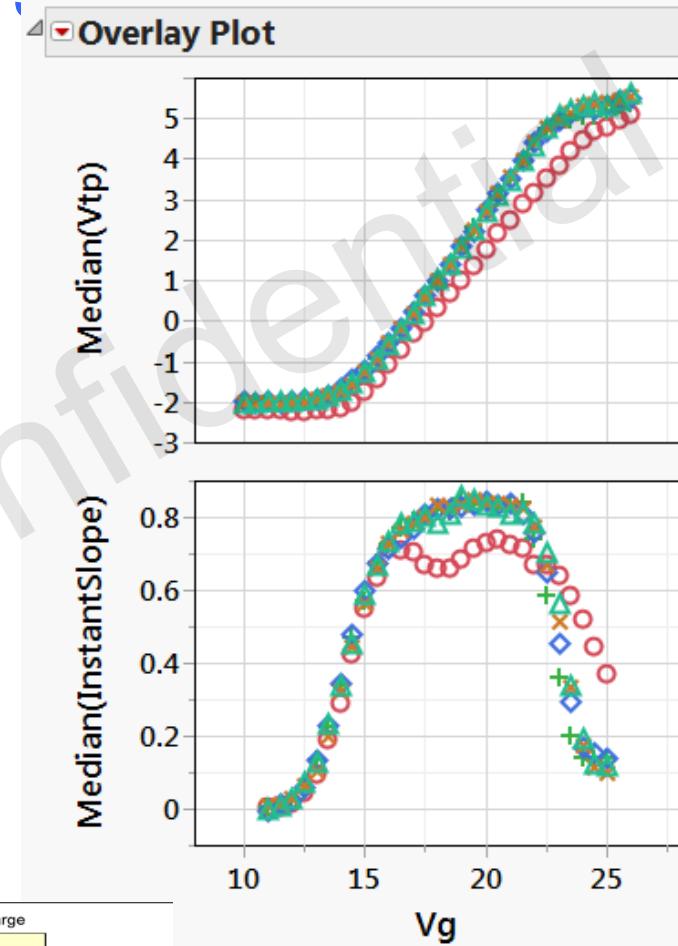


- E-field across TunOX and IPD are modulated due to the curvature effect.
- E-field at inner side of the interface increases and it decreases at the outer interface.
- Impact on P/E window is small.
Inner(tuox) vs. Inner (IPD) compete for V_{tp_sat} .
Outer(tuox) vs. outer (IPD) compete for V_{te_sat} .
- Impact on V_{gVt} and V_{wVt} is asymmetric.
 V_{gVt} has stronger diameter dependency. GCR and E-field effects are additive.

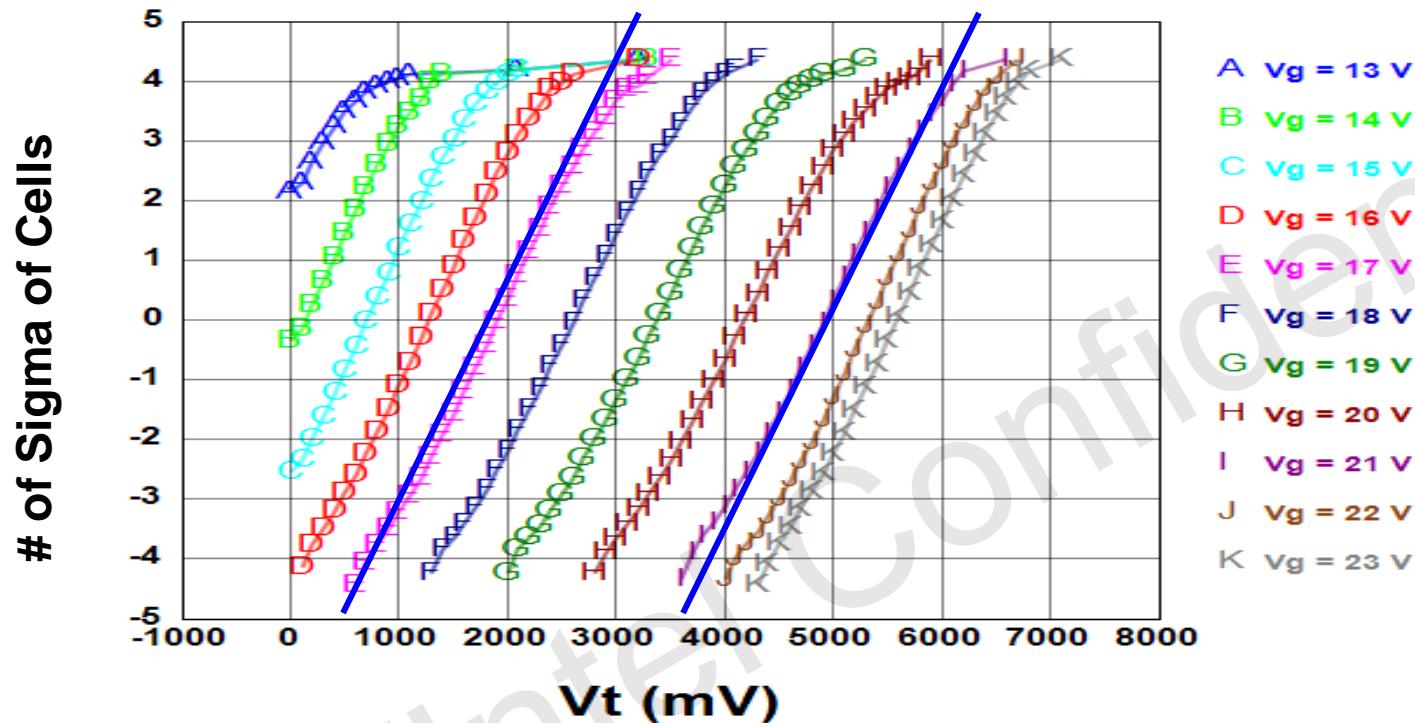


Programming Conditioning

- 3D NAND suffers from low GCR and hence low P/E window
- We can use all the help we can to improve P/E window
- Trapping some charge in the middle Nitride (by subjecting Cell to Program sat condition) of IPD result in reduced IPD leakage and thus improved Programming and V_{tp}_sat in subsequent programming operations
- Pre-conditioning/FCC: using several targeted program/erase cycles at probe to “condition” the cell for better programming capability.
- It is critical that conditioning doesn’t go away during assembly bake or reliability bakes.



Array Programming Characteristics

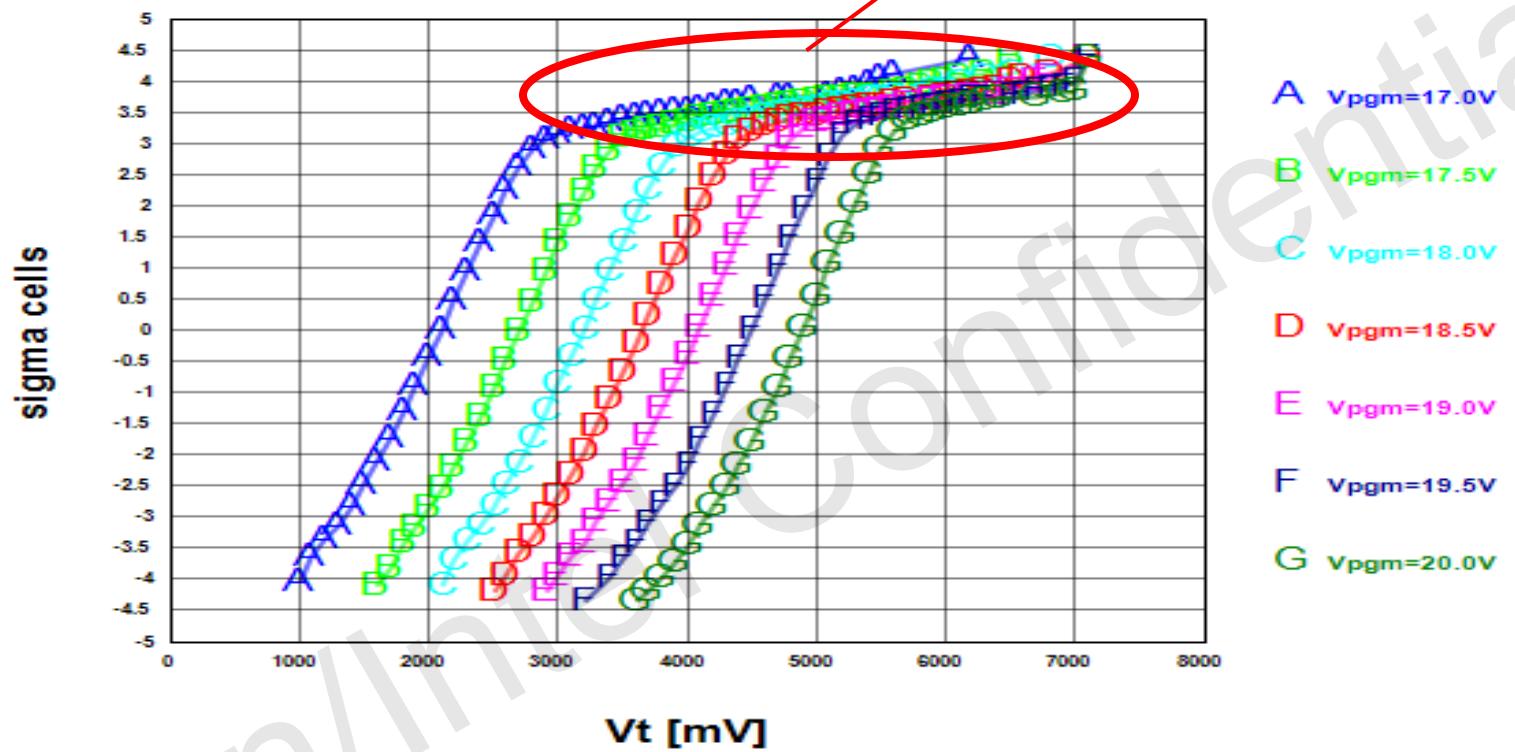


- Cells in array program to different V_t 's when programmed using a given gate voltage pulse
 - Fast cells are those that program to highest V_t , while slow cells are those that program to lowest V_t , for a given gate voltage
 - Slope of distribution is PVS (Program Vt Sigma)
- Fast cell and slow cell move by same amount for subsequent pulses, maintaining the difference in V_t 's reached for a given gate voltage

Program(Erase) Distributions

- PVS is one of the most important parameter for NAND cell. It sets the program time, Program disturb bit error rate, max voltage requirements, as well as dielectric scalability. It is very critical to minimize PVS
- Program/Erase Distribution are affected by several factors –
 1. Native V_t variations in the cell – random doping/trap variations
 2. Tunneling current variations from cell to cell due to - Tunnel barrier variations, Si/Poly1 asperities, etc.
 3. Cell Coupling Ratio variations – variation in W, L, Recess, Tox , T_{ono} , etc.
 4. Others - Bird's beak, systematic effects from design
- Empirically we see that the PVS gets worse with cell scaling – almost as $1/\sqrt{W*L}$. Here item #1 and #2 would be expected to scale as $1/\sqrt{W*L} \rightarrow 1/k$.
- Our experience suggests that, item #3 and #4 can be fixed/improved to the level that they are not fundamental limiters.

EOP – Extrinsic Over Programming



- In addition to the intrinsic PVS, 100s also showed extrinsic Over Programming tail
- Was attributed to the tunnel-oxide weakness (electrical) requiring optimization of TB dep/cut

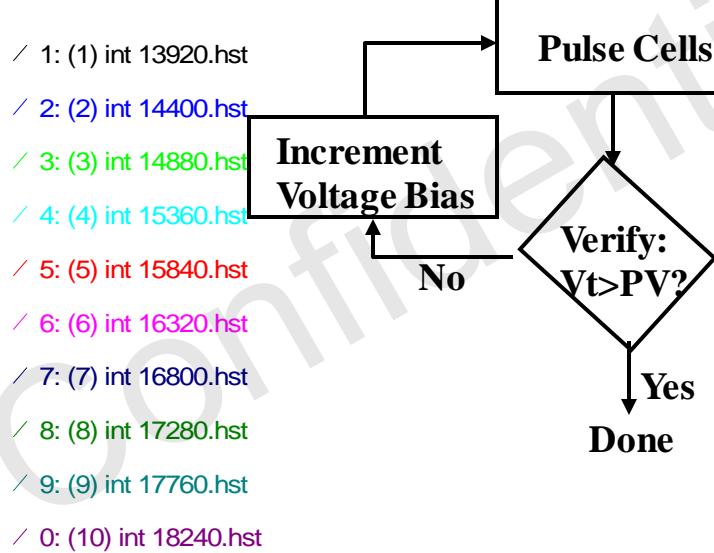
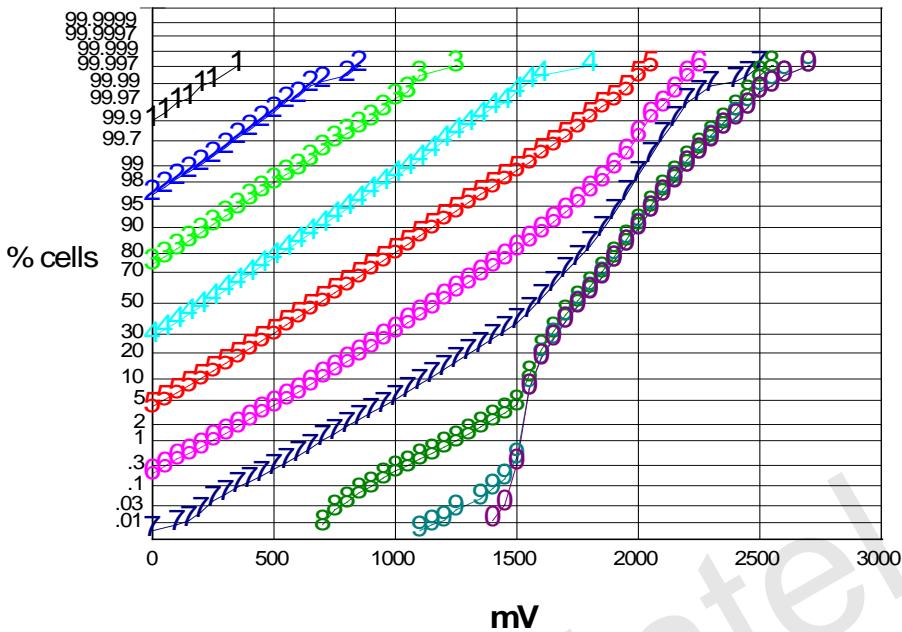
Programming Algorithm

- Since different cells in the array program to different V_t 's when programmed using a given gate voltage pulse, we need to use some sort of algorithm to be able to program the cells to tight program distribution
- Programming (Placement) algorithm:
 1. Start with a low programming gate voltage
 2. Give a programming gate pulse (~20us)
 3. Check the cell V_t (Verify operation).
 - If cell $V_t >$ Target V_t , end programming loop
 - If the cell $V_t <$ Target V_t , increment the gate voltage (~0.5-1.0V).
 - Repeat steps 2-3.
- Starting gate voltage needs to be low enough to ensure the cell does not over program (dictated by lowest V_g-V_{tx}) in the first pulse
- Final gate voltage is dictated by the cell that has the highest V_g-V_{tx}
- Total number of programming pulses, and hence the programming speed is dictated by the spread in the V_g-V_{tx} (\Rightarrow PVS)

Program- Pulse/Verify

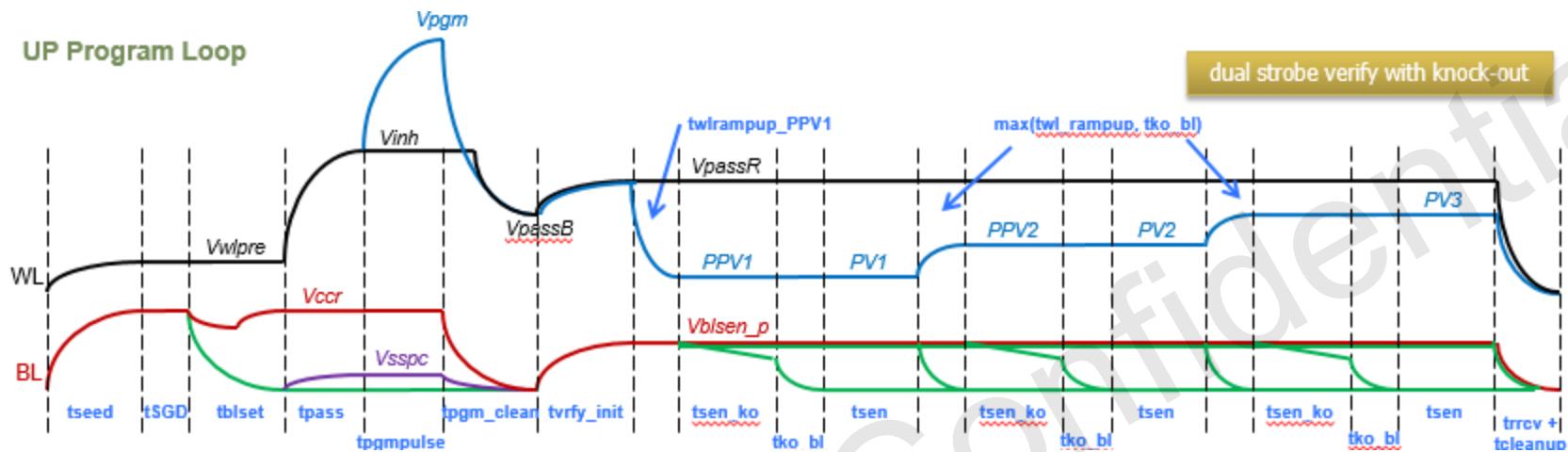
Gaussian

NAND Program Algorithm



- Each cell in 16kB Page Buffer is twiddled off when cell successfully Verified. Example above shows cell being programmed using ~0.5V gate step starting with a gate voltage of 13.9V.
 - Fastest to program cell crosses program verify (~1.5V) at $V_g=15.4V$
 - Slowest to program cell crosses program verify (~1.5V) at $V_g=18.2V$
 - Ideally this should give a program V_t distribution width that is 0.5V wide. However, the actual distribution is >1V (due to Program Noise (see later))
 - Without a pulse/verify algorithm the distribution would have been 3-4V wide

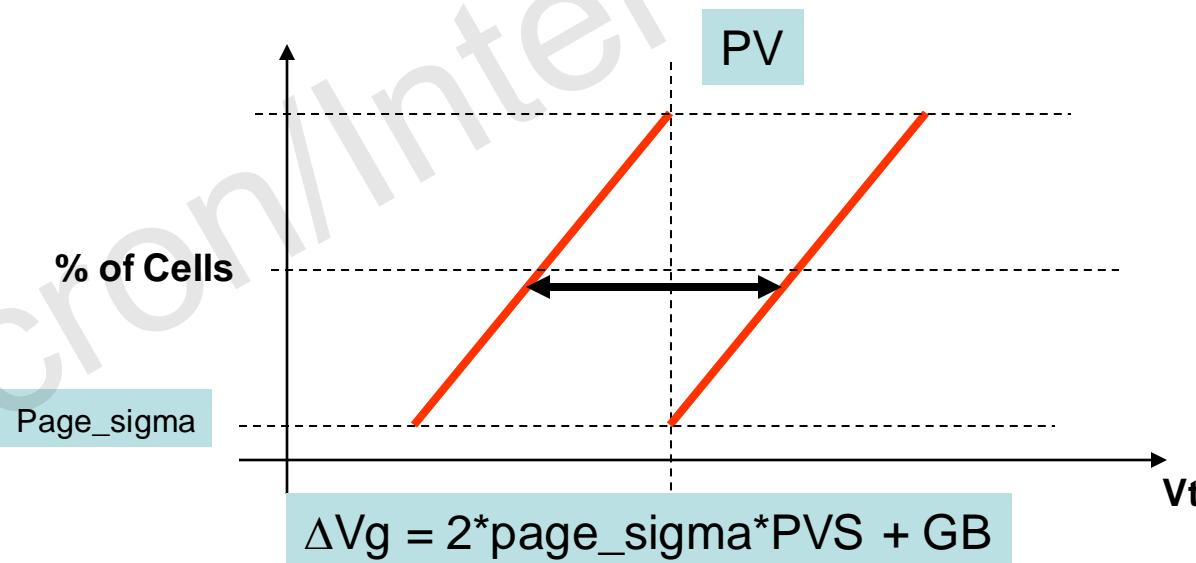
Simplified Programming waveform



- Total Program time (t_{PROG}) is # of pulses * loop-time
- We discussed number of pulses earlier (dependent on PVS, gate-step, cfbyte)
- The Loop time is made up of the above components, and each is carefully optimized.

Program Speed: Page Pgm

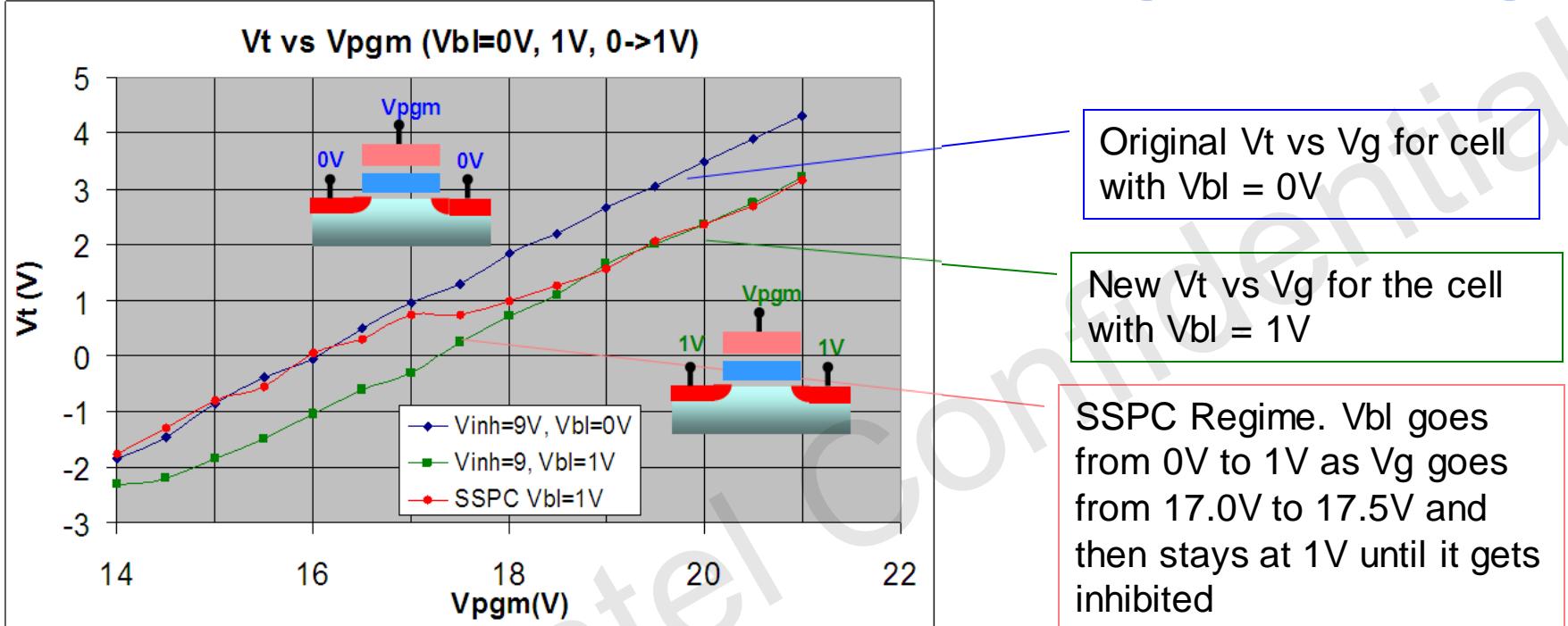
- Program Time
 - The gate step will determine the program V_t distribution that can be achieved. Smaller the gate step tighter the Program V_t distribution
 - Increased PVS leads to a smaller V_{gVt_min}, and larger V_{gVt_max}
 - As a result, the V_{g_start} will have to be lower, and V_{g_end} higher leading to increased number of pulses and increased program time
 - $V_{g_end} - V_{g_start} = 2 * \# \sigma * PVS + \text{Guard-band} (PV_{max} - PV_{min})$
 - $\# \text{ steps} = \Delta V_g / \text{Gate_Step}$



SSPC – Selective Slow Programming Convergence

- If we step up the program gate voltage by ΔV_g and give a programming pulse, the cell V_t will also increase by the same amount (ΔV_g)
- SSPC is a scheme that will cause the cell to move by less than ΔV_g even though the Gate voltage increased by ΔV_g .
- This is done by biasing the S/D/Channel (Bitline) of the cell of the cell to small voltage ~0.5V – 1.0V.
- Since the Tunneling current (Programming current) is actually a function of the FG to Channel field, even though the V_g is increased by ΔV_g the actual channel to control gate voltage is increasing only by $(\Delta V_g - \Delta V_{bl})$ and as a result the increased V_g doesn't result in the corresponding increased current -> slow down of programming.
- Good way to achieve program distributions tighter than the gate step. Helps with programming speed as the gate step can be large without compromising the program placement width.

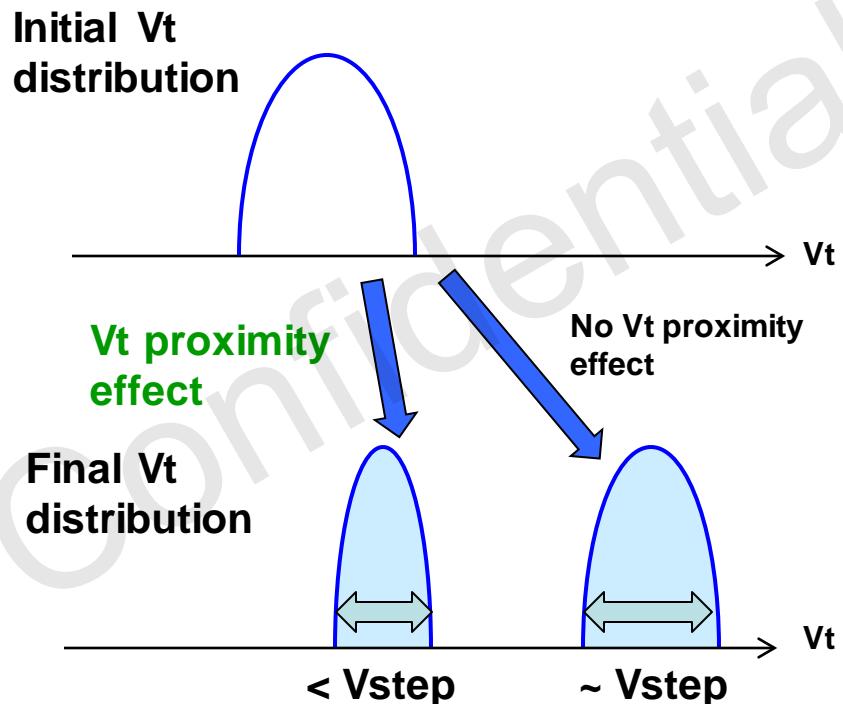
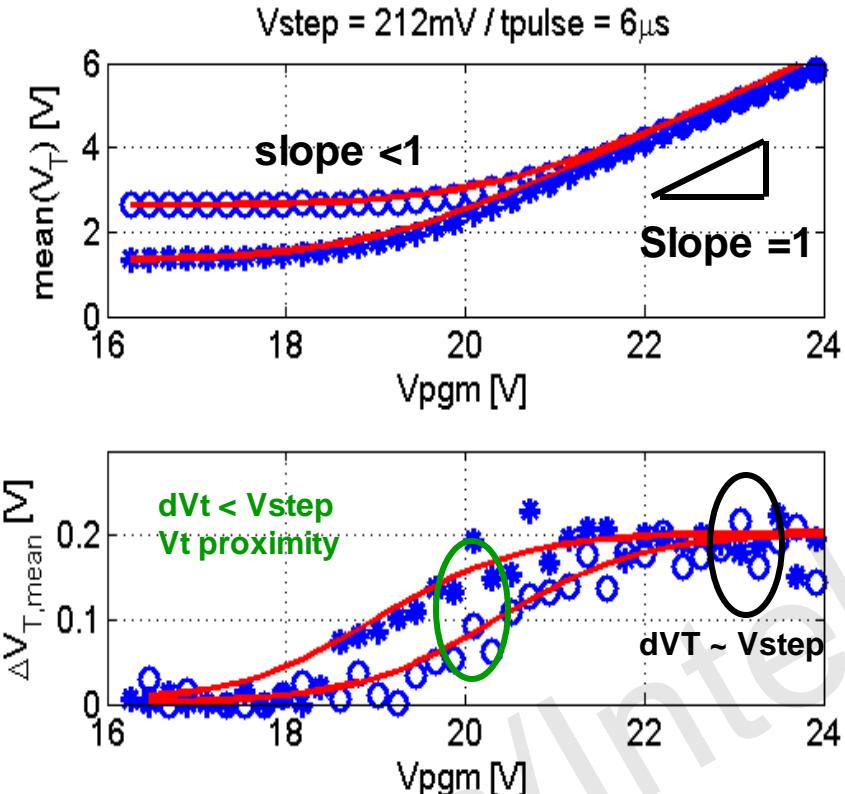
SSPC - Selective Slow Programming



- As the cell reaches close to the target V_t , (as determined by passing the pre-verify), SSPC voltage is enabled: Vbl goes from 0V to the SSPC voltage ($\sim 0.5\text{-}1.0\text{V}$).
- This slows down the programming of the cell for the next couple of pulses, allowing a finer placement of the cell V_t .

$$\Delta V_t / \Delta V_{bl} = (1 - CCR) / GCR \sim (1 + 2 * \text{WL-FG Coupling}) \rightarrow \text{"Magic Ratio"}$$

Vt proximity effect

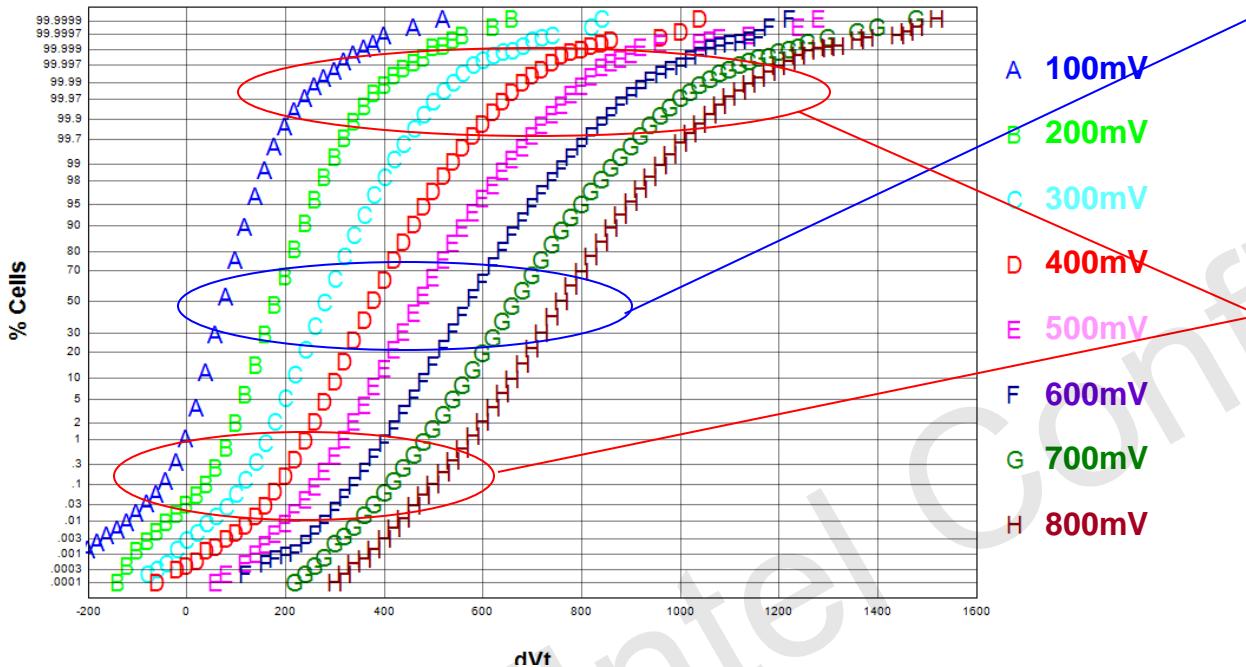


- If the program gate voltage starts low ($V_{g_start} < V_{t_init} + V_{gVt}$) initially cell moves by $\Delta V_t < \Delta V_g$ (Pgm slope < 1). Once $V_g > V_t + V_{gVt}$, the cell Vt moves at the same rate as ΔV_g (converges to Program Slope of ~ 1).
- If the final Vt is not far from initial Vt, a tighter state width is achieved.

Programming Noise/Fluctuation

Gaussian

dVt vs MC on L61



- When programming a large number of cells using a staircase gate algorithm (Gate step = ΔV_g), a typical cell V_t will move by $\sim \Delta V_g$ from step to step. However, we find that all cells don't move by this amount and there is a distribution of Delta_Vts – A normal distribution followed by a tail
 - Normal Distribution comes from number fluctuation mostly dictated by cell capacitance that will dictate the number of electron/program step
 - Tail distribution comes from erratic over-programming (random events of enhanced tunneling). This part tends to increase with cycling as oxide degrades

Programming Noise/Fluctuation

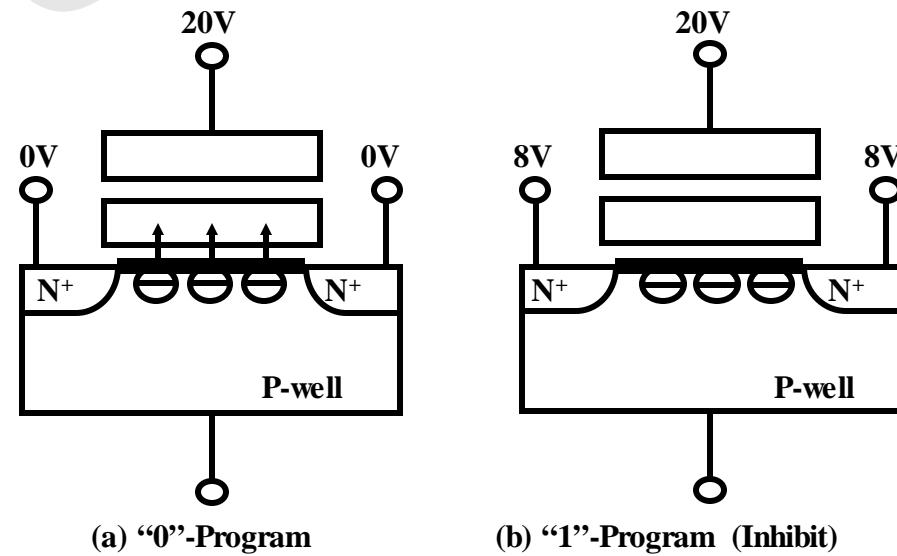
Tech Node	W	L	Tox	Tipd	Cipd	# Electrons/V	Δ_Vg	# Electron	Prog Noise
40s	55	72	75	135	4.39E-17	275	0.400	110	38
50s	40	50	75	130	2.54E-17	159	0.400	64	50
60s	30	35	75	120	1.71E-17	107	0.400	43	61
70s	28	24	75	110	9.67E-18	60	0.400	24	81
80s	20	20	60	85	5.31E-18	33	0.400	13	110
100s	251	15	70	130	2.96E-17	185	0.400	74	46

- A simple 1st order calculation can help us compute the Programming noise which is in reasonable agreement with data
 - During every programming step we put certain number of electrons on to the floating gate
 - Since $\Delta V_t = \Delta V_g$, and $\Delta Q_{fg} = \Delta V_t * C_{ipd}$, we can compute the amount of charge added.
 - Dividing this by charge of an electron we can calculate the number of electrons added to the FG on the average every pulse
 - The standard deviation is $\text{sqrt}(\text{mean}) \rightarrow \Delta V_t_{\text{sigma}}$
- This is an approximate calculation. More exact calculations incorporating FN equation can be done as well.
- Higher cell capacitance of 100s cell leads to less Program noise

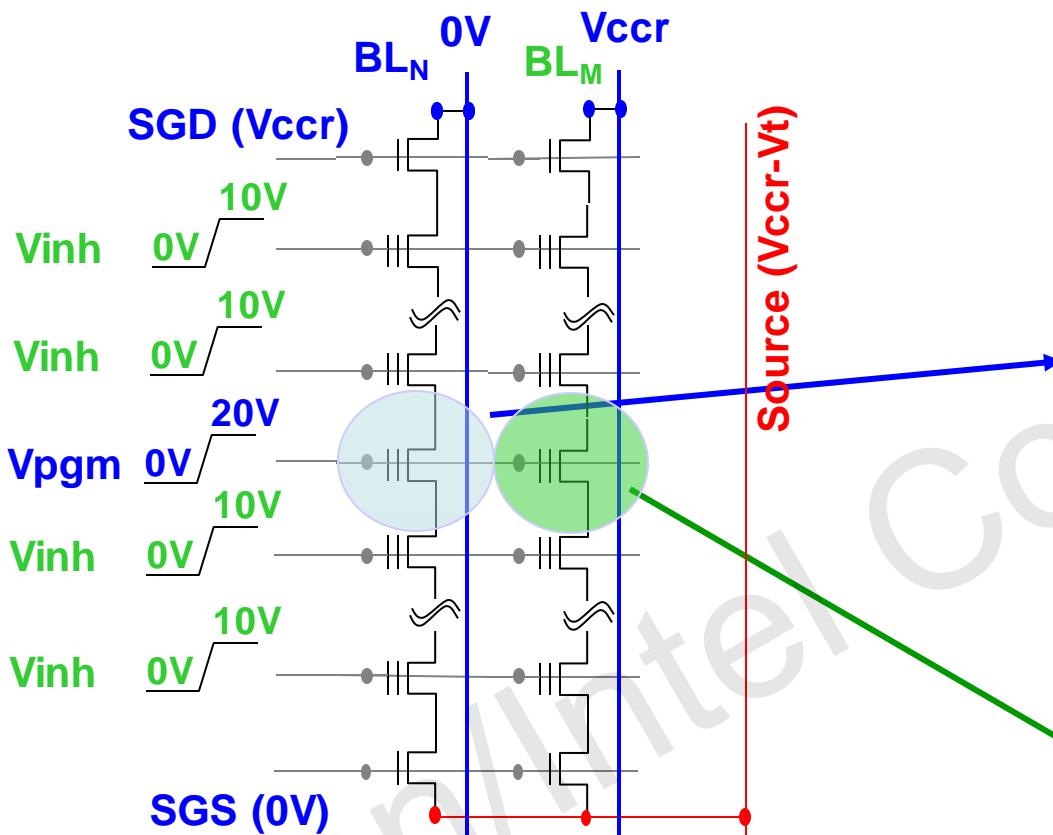
Inhibit

- We described the NAND Programming as FN tunneling from channel to the FG with the channel grounded and the WL high at 20V
- So, one question is how do we prevent (“inhibit”) the other cells on the WL from programming up as well, since we would not want all cells on the WL to get programmed up
- The way we achieve inhibiting programming on all the other cells on the selected WL is called the “Program Inhibit”

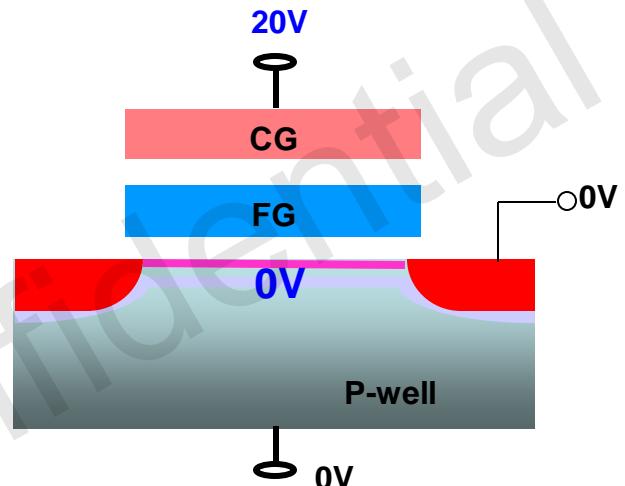
- Program inhibited on other cells on the same wordline by floating those string and letting the channel and S/D “boost” up to ~6-8V, thus reducing the field across the Tunox



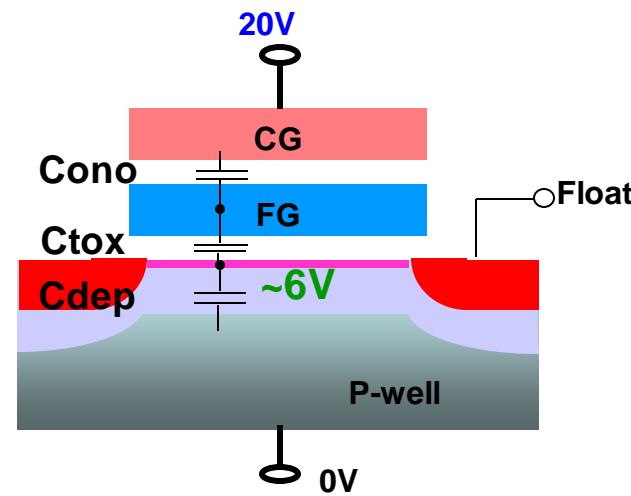
Program Inhibit (Self-Boosting)



Cell to be Programmed

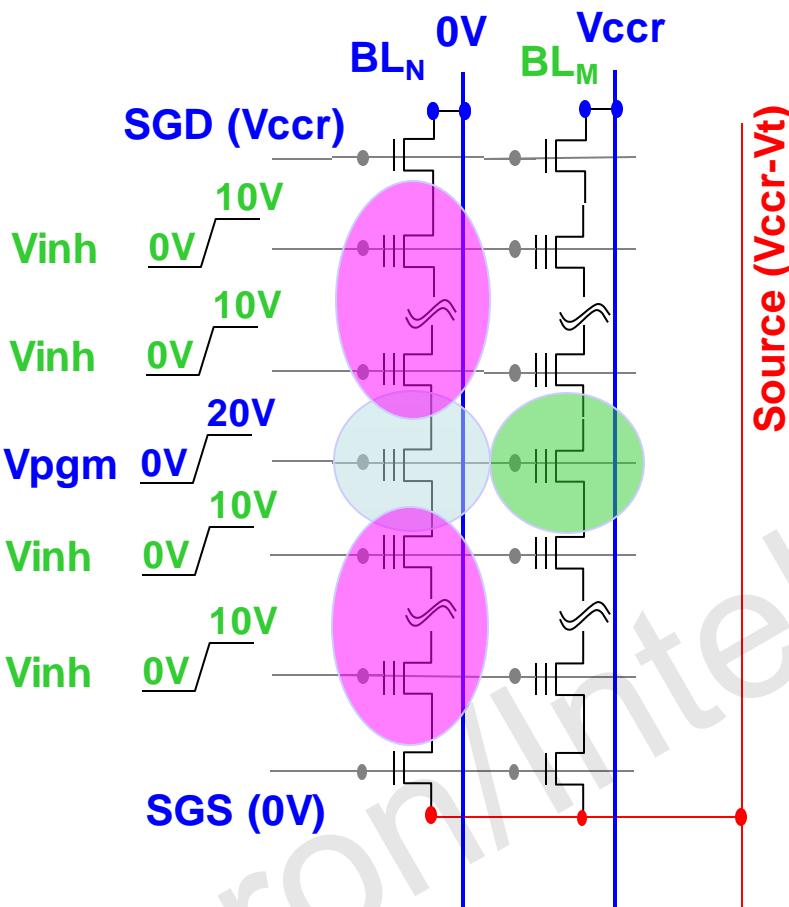


Cell to be inhibited



- Uses the capacitive coupling between the Control Gate and the Channel to boost the channel to high enough a voltage to inhibit programming

Program/Inhibit Disturb

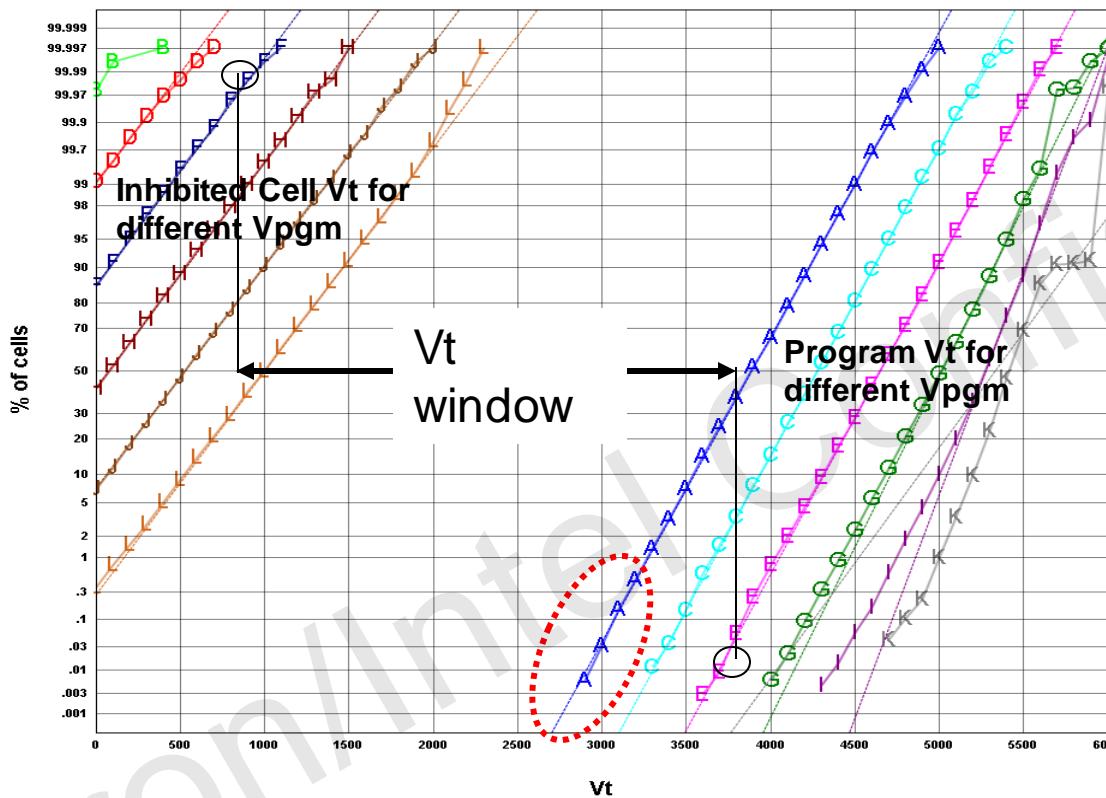


- There are two types of disturbances during programming:
 - Cells on the same WL see program disturb. If the Programming WL voltage is high and boosted channel voltage is not high enough these cells can get disturbed up
 - Cells on the selected bitline see Inhibit disturb. Here the channel is grounded and all the deselected WL are sitting at Inhibit gate voltage. If the Inhibit gate voltage is high, these cells can get disturbed up
- Program Disturb happens when the cells on that WL are programmed
- Inhibit disturb happens every time any other WL in that block is programmed

Program Disturb Window

Gaussian

L52A WL10 Cyc=1 Temp=25c Vinh=5910mV



A Page: 44 Voltage: 21400 pulse: 10, 30-99.9%, M=3.89e+03, S=263

B Page: 45 Voltage: 21400 pulse: 10

C Page: 44 Voltage: 21800 pulse: 11, 30-99.9%, M=4.27e+03, S=261

D Page: 45 Voltage: 21800 pulse: 11, 30-99.9%, M=861, S=366

E Page: 44 Voltage: 22200 pulse: 12, 30-99.9%, M=4.65e+03, S=257

F Page: 45 Voltage: 22200 pulse: 12, 30-99.9%, M=381, S=354

G Page: 44 Voltage: 22600 pulse: 13, 30-99.9%, M=5.01e+03, S=236

H Page: 45 Voltage: 22600 pulse: 13, 30-99.9%, M=71.9, S=357

I Page: 44 Voltage: 23000 pulse: 14, 30-99.9%, M=5.27e+03, S=179

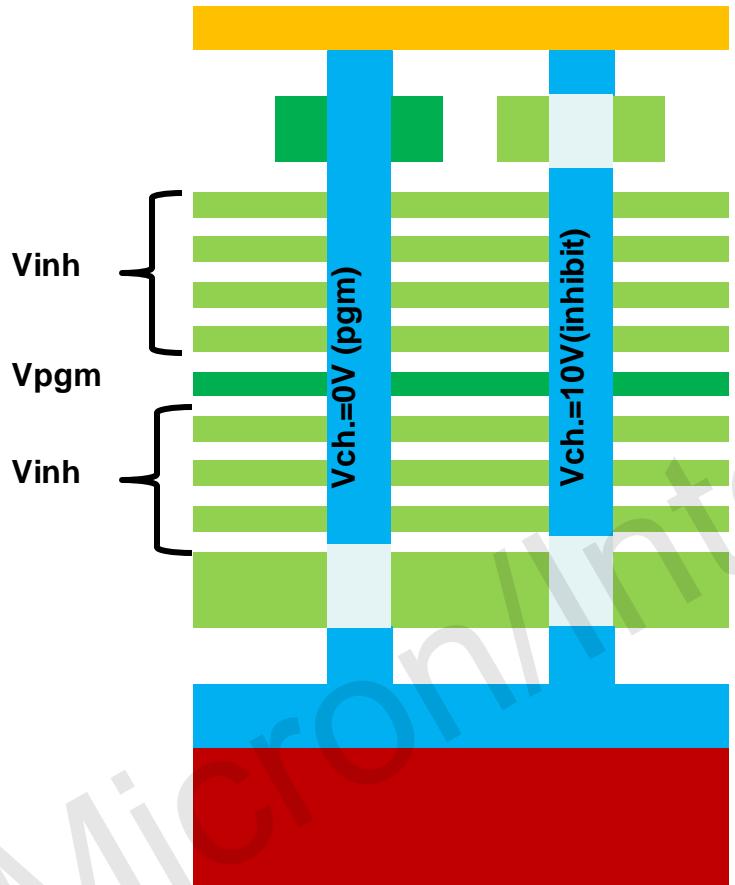
J Page: 45 Voltage: 23000 pulse: 14, 30-99.9%, M=525, S=362

K Page: 44 Voltage: 23400 pulse: 15, 30-99.9%, M=5.33e+03, S=348

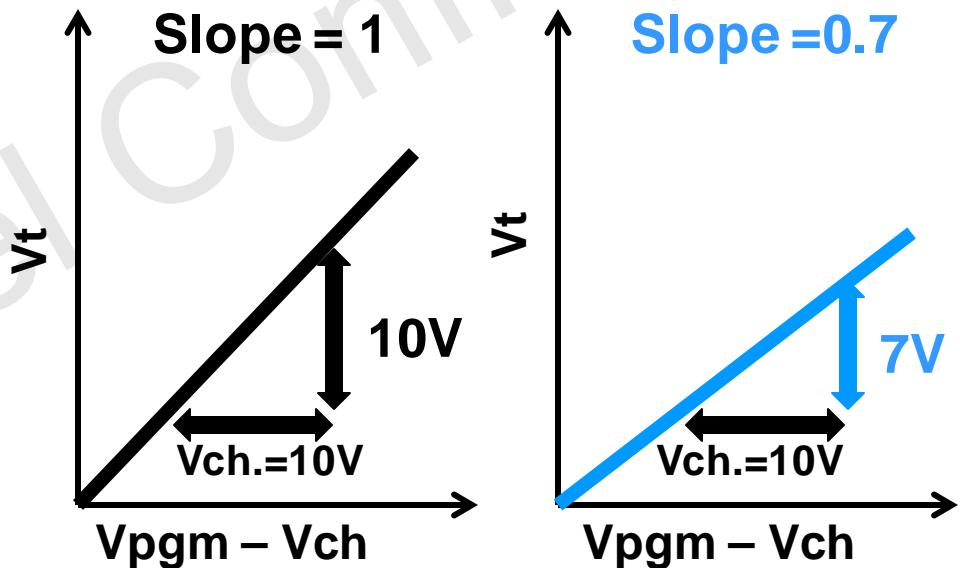
L Page: 45 Voltage: 23400 pulse: 15, 30-99.9%, M=998, S=360

- To guarantee good Program Disturb Window, we need tail to tail window > PV7
- Intrinsic Program Disturb window > PV7 – R1 + ~7*PVS + FG-FG GB

Programming slope requirement for program disturb window



Program disturb V_t window



Pgm slope <1 makes pgm disturb window more challenging.

Boosting Requirements

- If Vinh is too low, we do not get enough boosting and we have program disturb fail
 - If Vinh is too high then we can have Inhibit disturb
 - The cells on the selected strings see the Inhibit disturb for the duration of 31 programming
 - For the high Vinh, in addition to the Inhibit disturb, we can also have high inhibit voltage induced disturbs – hot-e disturb, GIDL disturb, punchthrough disturb, etc.
- So, need to have adequate window between these two where we can avoid disturb

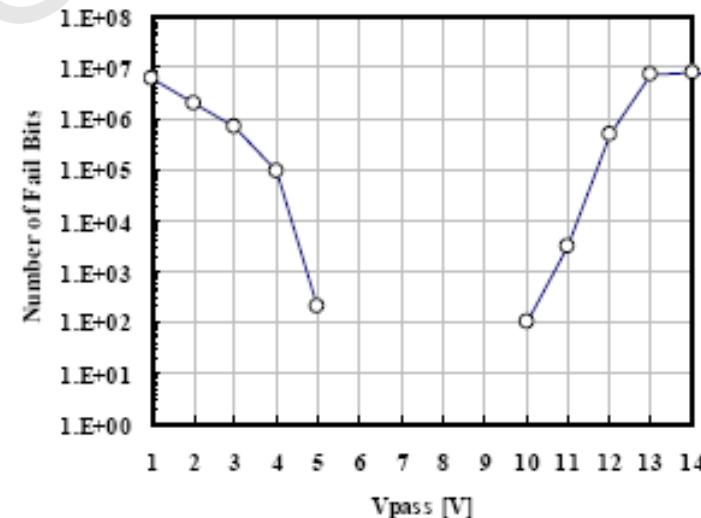
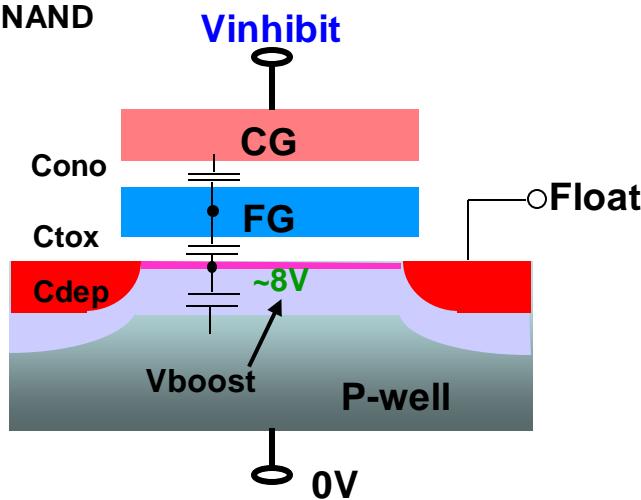


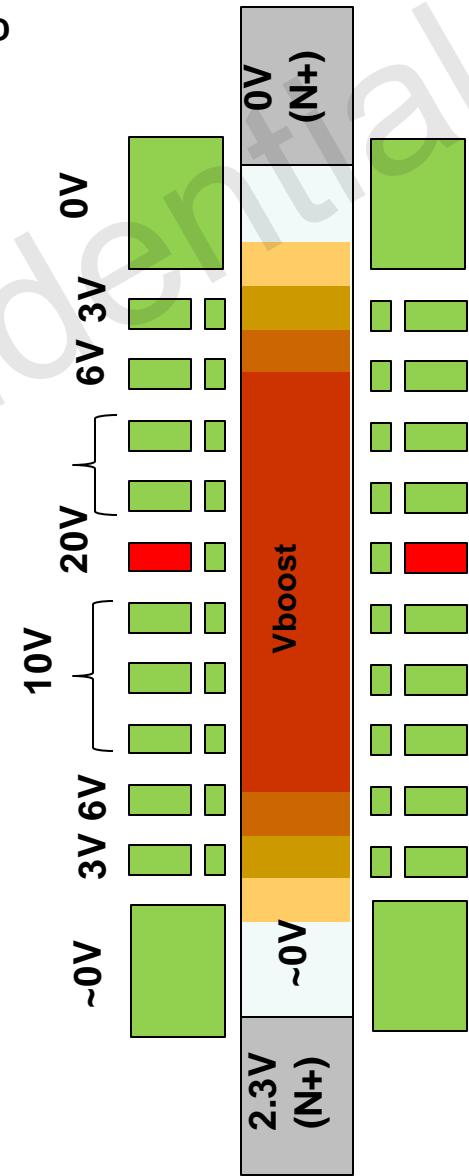
Fig. 10. Program disturbance of the NAND cell arrays with 70nm design rules as a function of the pass voltage.

2D vs 3D Program Disturb (PD)

2D NAND



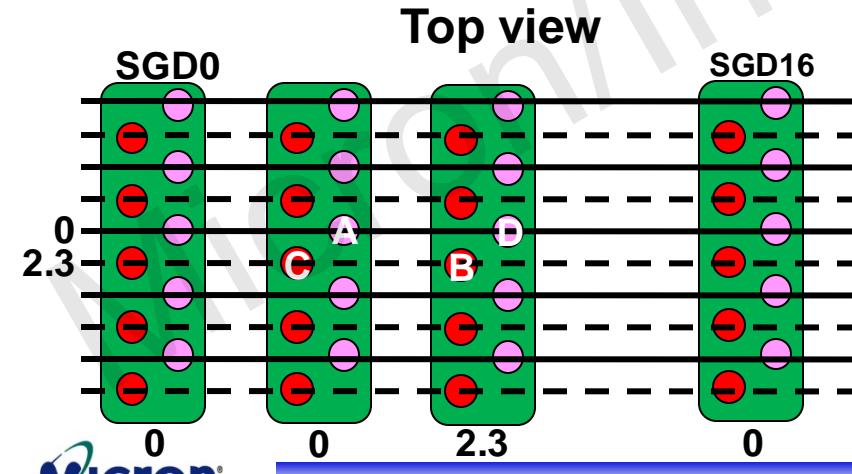
3D NAND



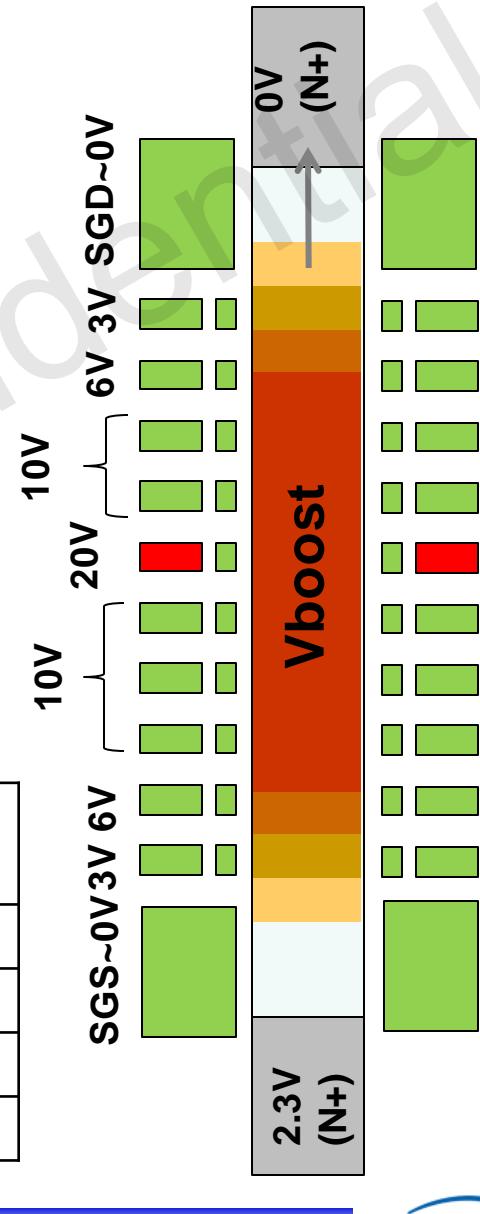
- Vboost created by capacitive divider – Vinh “fights” the depletion capacitance of the grounded P-well
- Coupling must also fight leakage mechanisms discharging Vboost
- 3D with Floating pillar has an advantage – no grounded P-well to fight → Boost ratio should be ~1
- But leakage still a risk, especially because of polysilicon pillar (grain boundaries and other defects)

Inhibit and boost

- No grounded P-well to fight: higher boost coupling ratio
- Leakage through the SGS or SGD is a risk, so blocking schemes at block edge are important
- Each WL split among 16 sub-blocks (16 SGDs)
 - 1x conventional inhibit (high SGD, high BL)
 - 15x new inhibit (low SGD, low BL)
- “A” especially sensitive to SGD leakage, because BL=0
- More efficient boost can lead to hot-e issues: on inhibited subblocks, high electric fields can be present, leading to h+/e- pairs and hot-e injection
- Careful optimization of Vpass voltage and inhibit scheme needs to be achieved to prevent this issue

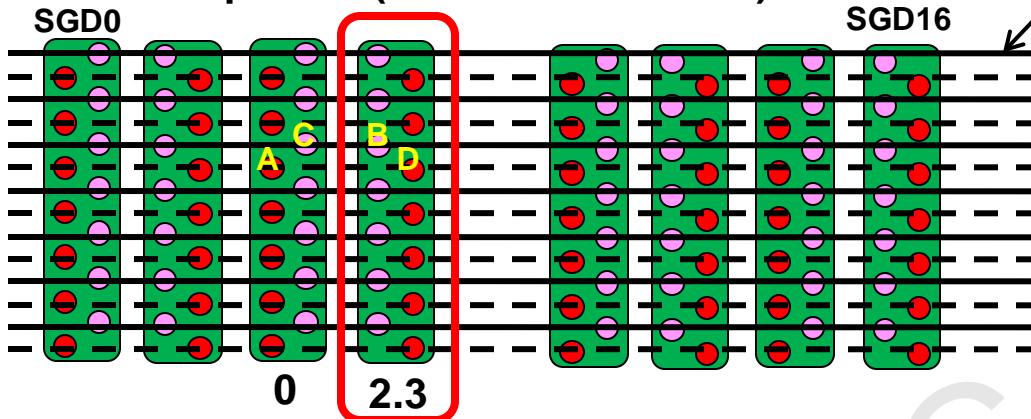


Pillar	Inhibit condition
A	New VSGD=0
B	Conventional
C	Stronger
D	Programmed



3D Program Disturb Biases

P-Block top View (14 Bit lines shown)



Bit Lines

2D case

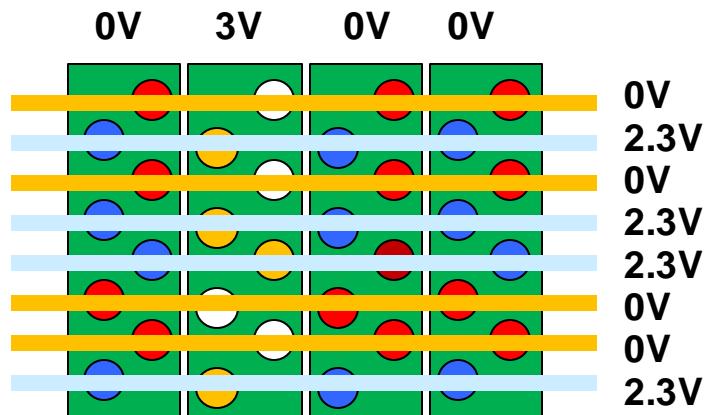
stress condition	WL	SGD	BL	pillar
A	vpgm	0	0	VboostA
B	vpgm	2.3	2.3	VboostB
C	vpgm	0	2.3	VboostC
D (programmed)	vpgm	2.3	0	0

Selected sub block
D=selected cell

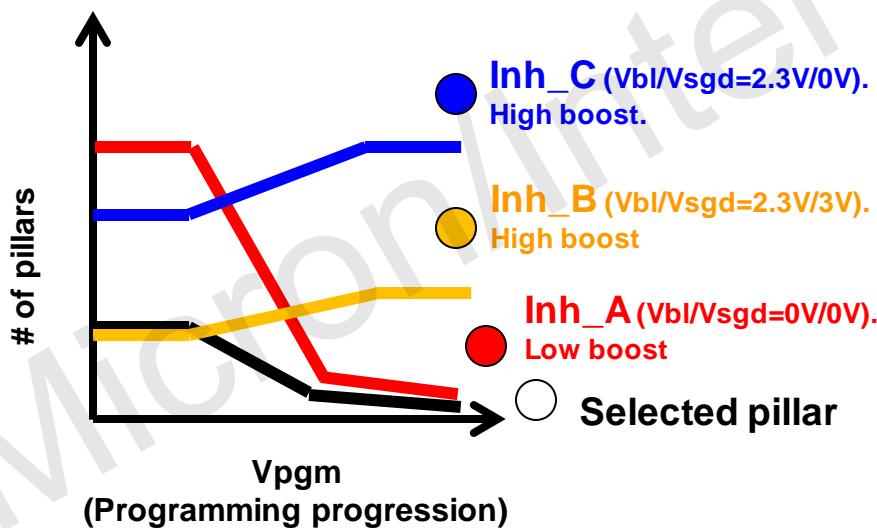
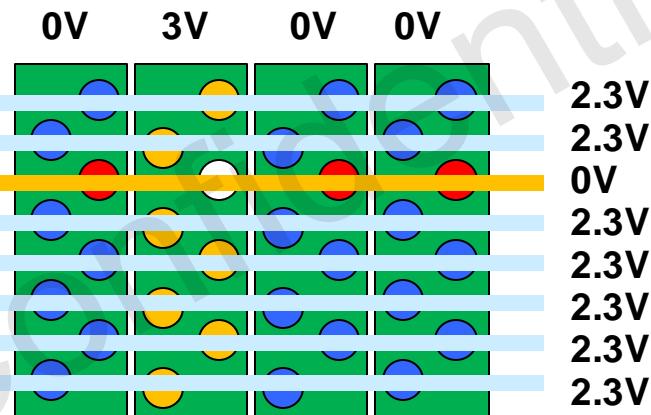
- 3D advantage is good boost ratio (surround gate and floating pillar)
- But 3D must inhibit cells on selected WL and deselected SGD (top, cases "A" & "C"). "A" worse due to the low BL voltage.
- Condition "C" expected to have good boost (low SGD leakage), but has higher requirement because final programming pulses are higher voltage

Boosting Pattern Progression

Beginning of programming
($V_{pgm}=\text{low}$)

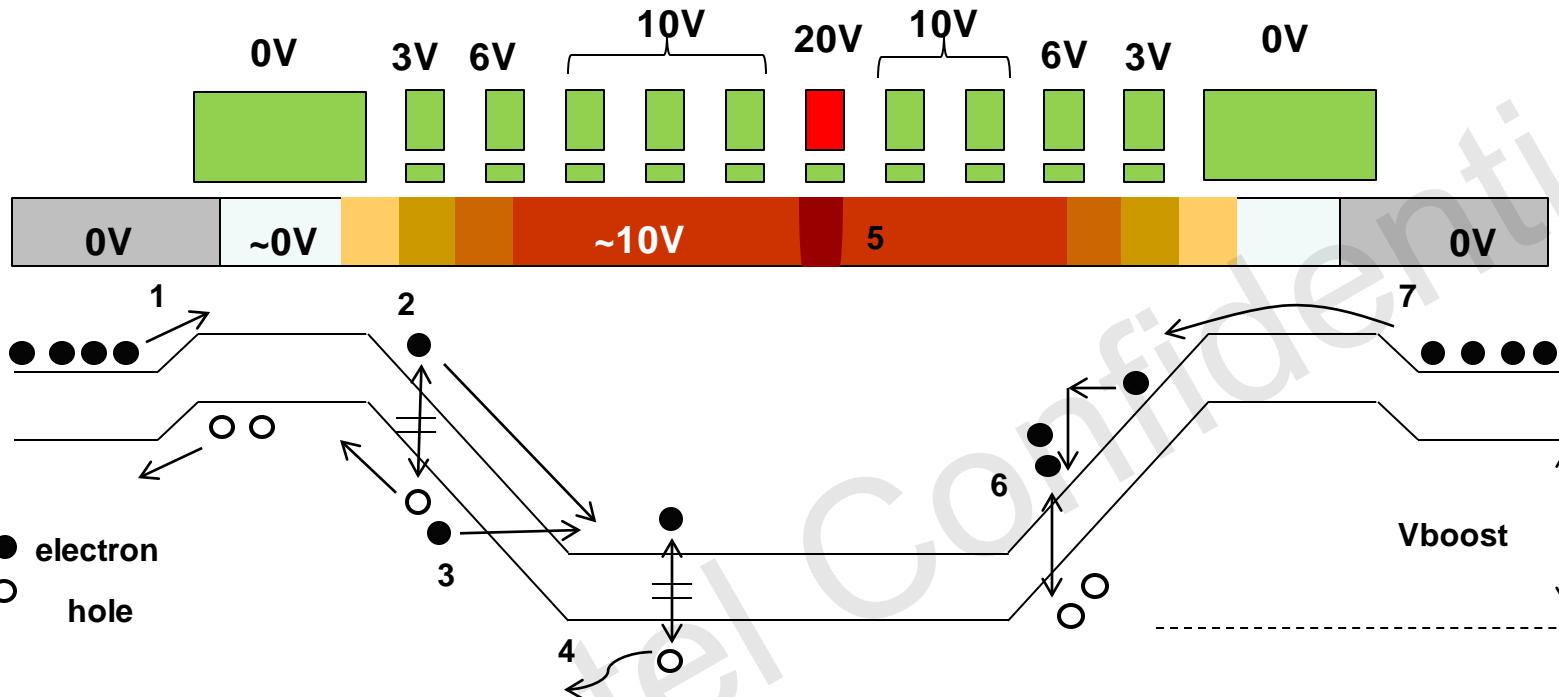


End of programming
($V_{pgm}=\text{high}$)



- As programming progresses, numbers of inh_A pillars decreases because more BLs are inhibited.
- Therefore boosting requirement is less for inh_A condition.
- Any of Inh_A,B and C can be a boosting limiter.

Boost loss mechanisms

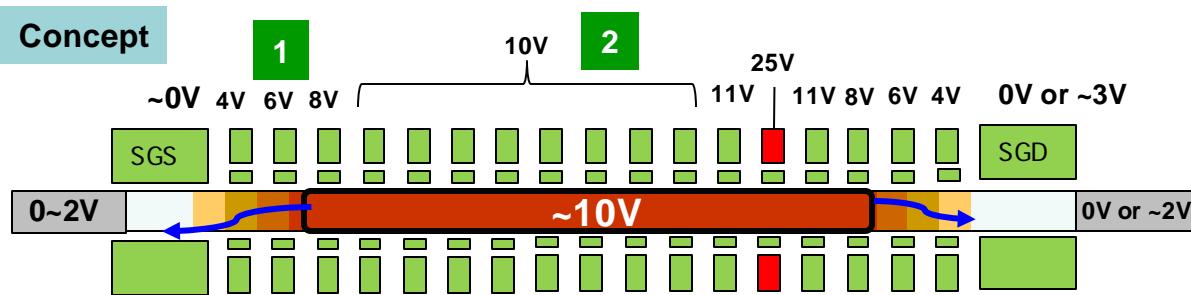


	Mechanism	Description
1	Bipolar gain	Holes under SG are injected into N+ then electrons are injected back.
2	e/h at edge	Thermally generated electrons are injected into channel due to E-field.
3	BTBT	Band-to-band tunneling at edge high E-field.
4	Hole diffusion	Thermally generated holes are diffused into edge before recombined with electrons..
5	Local boost hot e/h	Ch. Under Vpgm WL are locally boosted and creates hot e/h injection.
6	Impact ionization	Electrons are accelerated by edge E-field and create e/h pairs.
7	SG S/D leak	Source-drain leakage of SG transistor.

Electron injection into or hole loss from the boosting area cause boosting loss.

Self boost and Hot-e

Concept



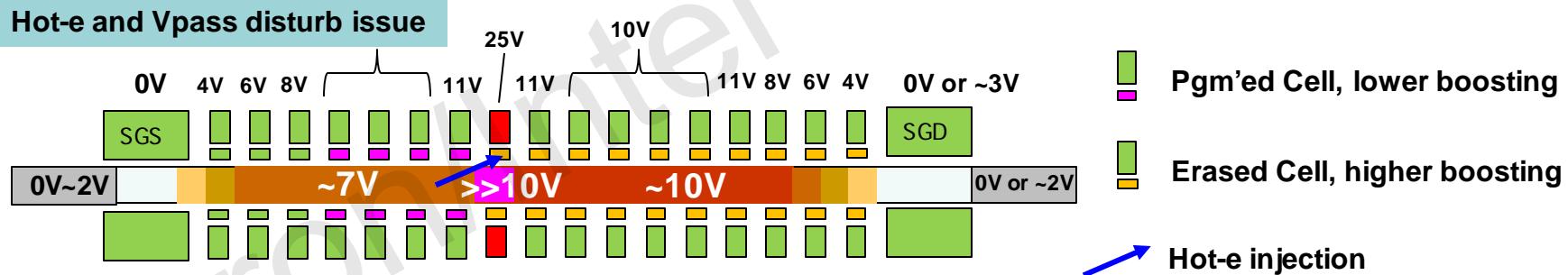
1. [Edge scheme] Edge WLs has gradual biasing to minimize the pillar E-field at the edge.

This reduces leakage at the edge, pulling down the boosting.

2. [Self-boost scheme] ~ All WLs are biased to ~10V, making the boosting area (capacitance) large.

This reduces the boosting loss with a given leakage at the edge. $dV = dQ/C$.

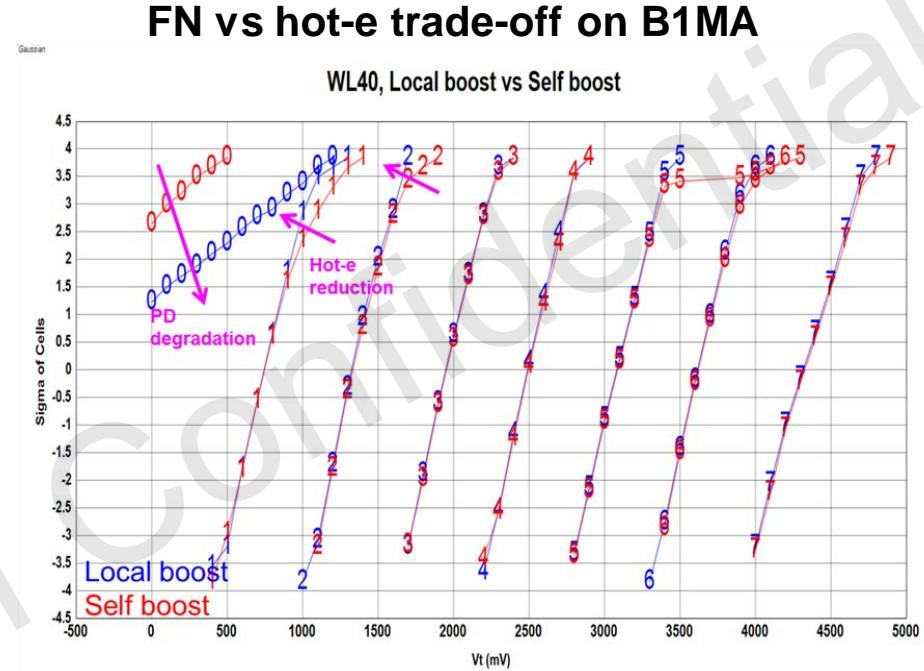
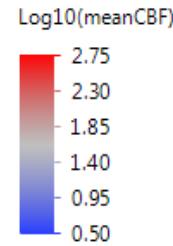
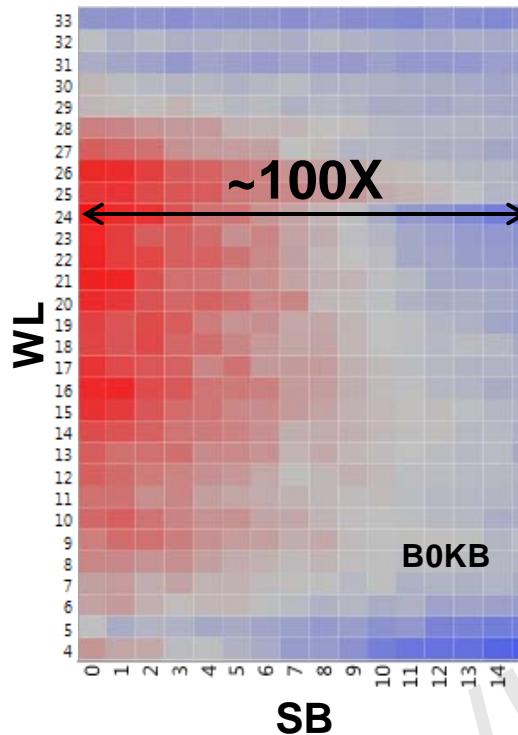
Hot-e and Vpass disturb issue



Hot-e injection disturb due to the lateral E-field.

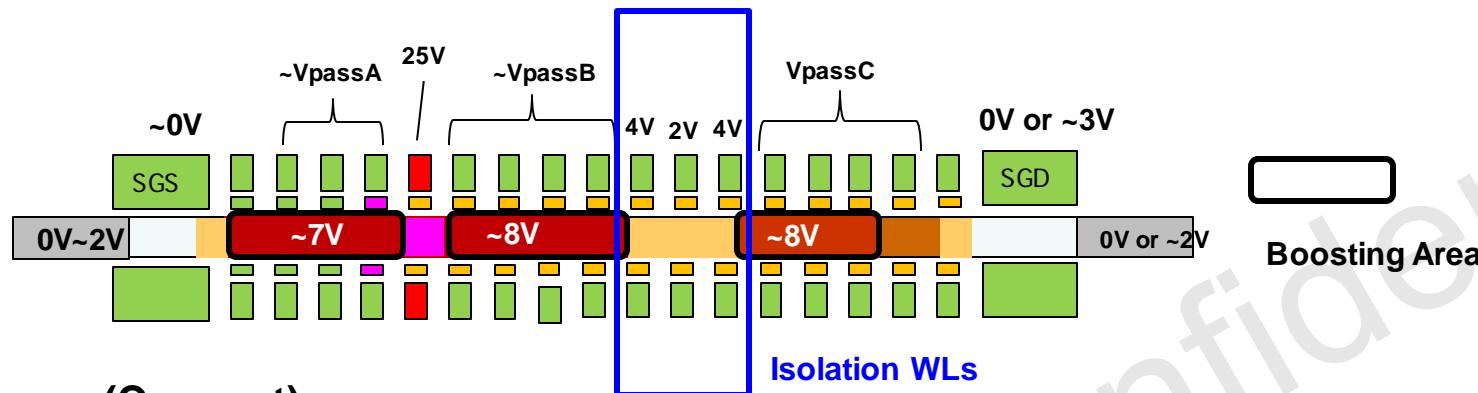
1. Data pattern is asymmetric, resulting in asymmetric boost potential.
2. Pillar under Vpgm is boosted higher due to WL-Pillar coupling.
→ The lateral E-field along the pillar create hot-e generation and injection to the cell.
3. Also, a high Vpass is applied to all WLs, degrading Vpass disturb.

PDOP: Hot-e impact on E2



- E2 is worse on center WLs (~WL24), and on the first subblock
- A clear trend is present with subblock (SB 15 on WL24 is ~100x better than SB0): the mechanism is impacting the subblocks that are inhibited after the setup
- The issue is caused by high Vpgm needed to bring center WLs to L7, leading to hot-e generation on already programmed subblocks
- Worst at 30C, uncycled, and it's improved by either cycling or high temperature

Local Self-boost



(Concept)

Segmenting the boosting area by putting **isolation WLs**.

Source side boosting scheme is shown above.

Based on the location of the iso. WLs,

Other variations such has Drain side boosting (iso. WLs are at SGS side) and Local self-boost (iso. WLs are at both sides) exist.

(Benefit)

Boost potential can be optimized by each boosting area depending on the data pattern and other circumstances. Vpass can be lowered in general, Improving Vpass disturb.

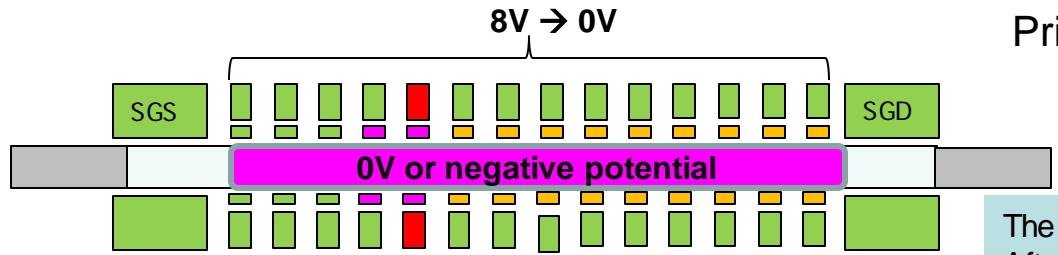
(Limitation)

Another lateral E-field can be generated around the iso. WLs.

The small boosting area (capacitance) can lose the potential easily.

Seeding

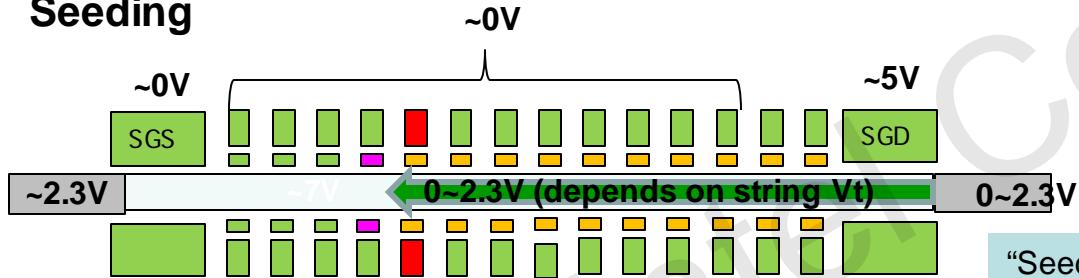
Before programming pulse
(after the previous pgm/verify loop)



“Seeding” sets the initial pillar potential
Prior to boosting.

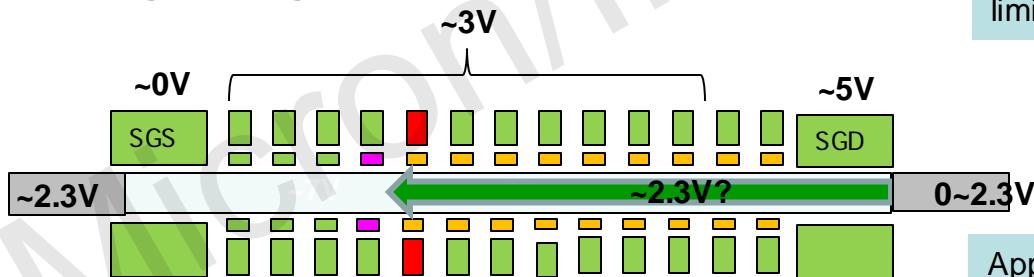
The pillar can be negative potential
After the previous pgm/PV, not desirable.

Seeding



“Seed” resets the pillar potential to 0V or
Higher depending on the string Vt.
Though all cells are the drain side is erased, the space Vt
limits the seed capability. (Similar to All WL EV)

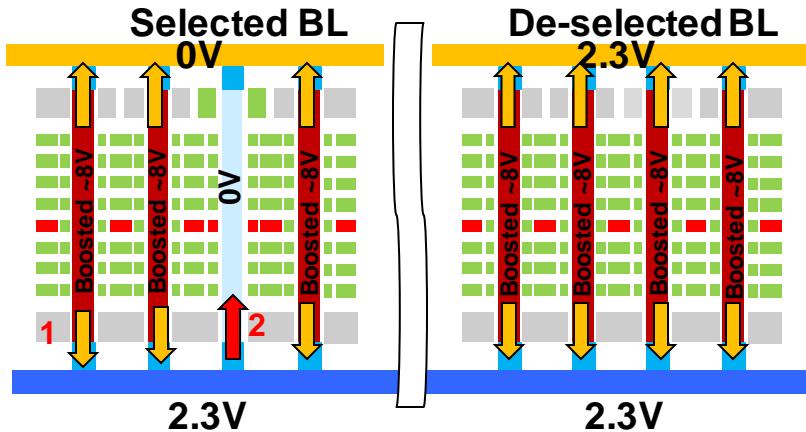
Seeding with high WL_seed



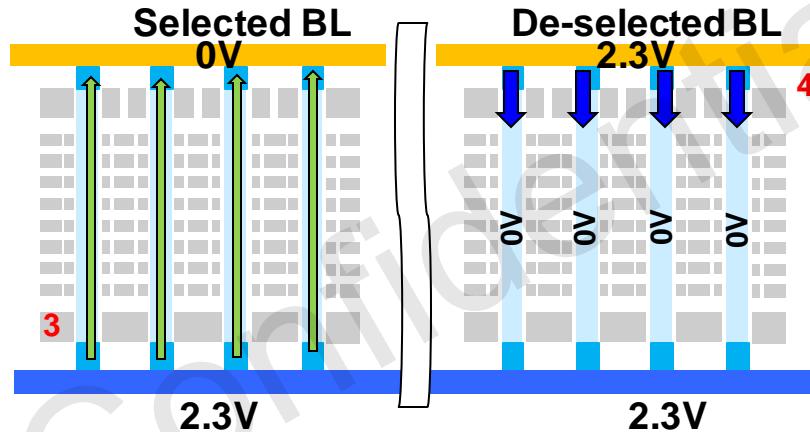
Applying positive WL bias for higher seeding often turns out
the net loss. Because you lose the boosting later.
+1V WL at seed provides;
≤1V seed gain vs.
~1V boost loss (WL-Pillar coupling is strong in 3D NAND)

SG leakage during programming

Selected Block



De-selected Block



1. SGS or SGD leakage pulling down boosting. Program disturb.
 2. SGS leakage pulling up the channel. Slow to programming.
 3. SGS leakage pulling up the channel then pulling up BL. Slow to programming.
 4. SGD leakage pulling down BL. Program disturb.

(*) SG leakage includes both source/drain punch through and GIDL.

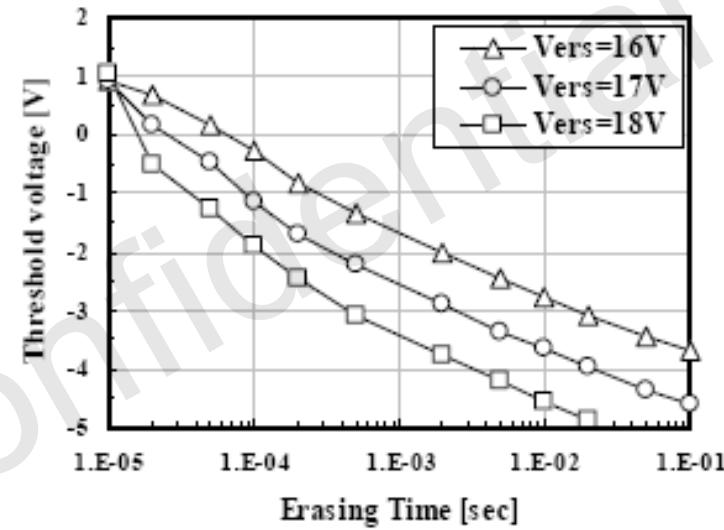
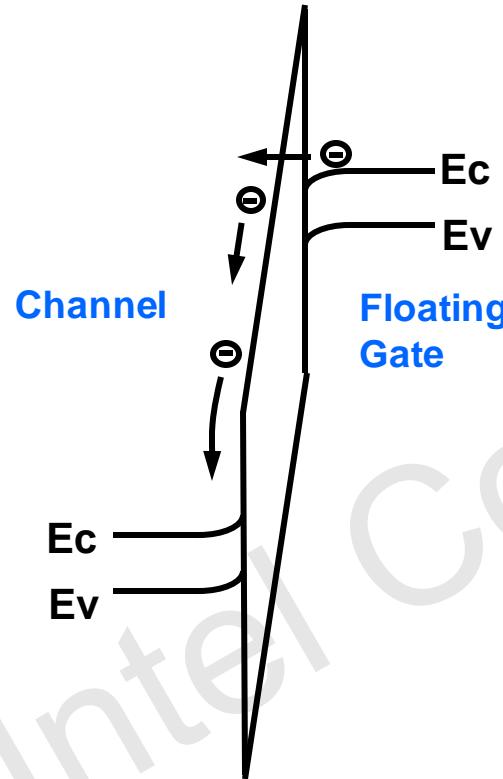
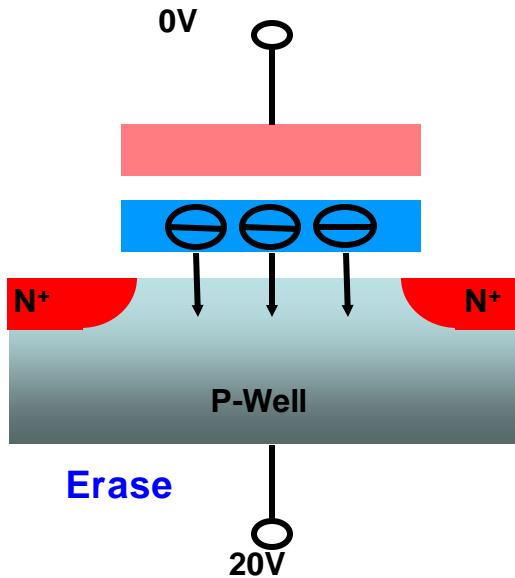
1,2 and 4 are common modes with 2D NAND.

3D NAND has increased risk due to more leakage caused by poly Si channel.

Erase

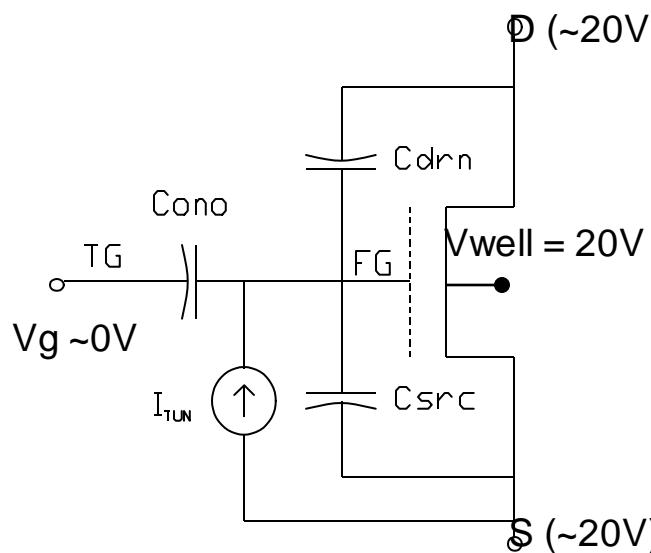
- FN Tunnel Erase
- Erase Inhibit
- Erase slope
- Erase Vt sigma
- 3D NAND Erase
- SG GIDL Requirements
- Erase Verify

NAND Flash Erase - FN Tunneling



- Tunnel Erase using Channel Erase with high positive voltage applied to the Well with gate grounded.
- Erase time ~1ms
- Erase Physics is essentially same as Program Physics, except now we tunnel from FG into the P-Well

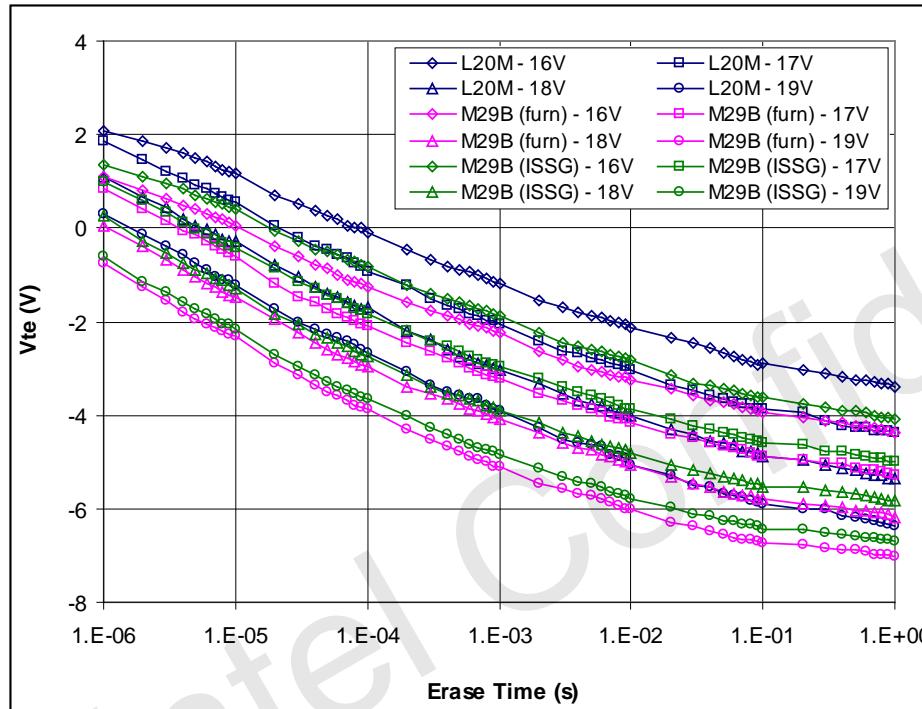
Erase V_T vs. Time



- Equiv. ckt shows electron tunneling off the Floating Gate by I_{TUN}
 - At the beginning of Erase, V_{FG} is low and I_{TUN} is large, so high tunnel current and erasing is fast
- $$V_{FG} = V_{CG} \cdot GCR + \frac{Q_{FG}}{C_{TOT}}$$
- $$V_T = \frac{V_{TFG}}{GCR} - \frac{Q_{FG}}{C_{ONO}}$$

- As number of electrons on FG decreases the FG voltage increases (by $\Delta Q_{fg}/C_{tot}$) decreasing I_{TUN} as given by the FN Equation
- Each $\sim 0.8V$ increase in V_{FG} drops I_{TUN} by $\sim 10x$
 - Erasing rate drops by 10x, 100x, etc as the floating gate gets charged by .8V, 1.6V, ... From the control gate perspective, the erase rate drops by 10x, 100x, etc as the Cell V_t decreases by $\sim 1.4V$, $2.8V$, ...
 - Makes Erase V_t vary as a "log time"

Erase

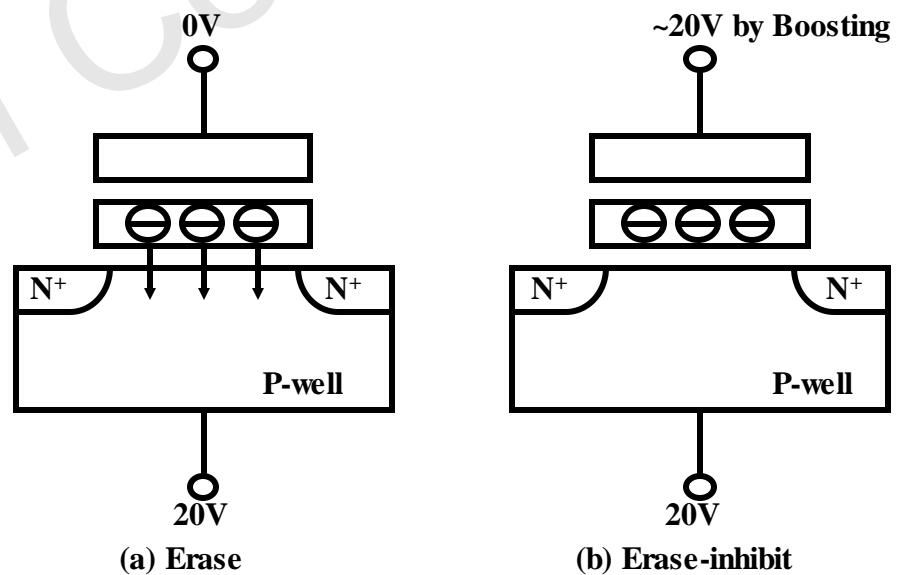


- Log Time behavior of the Erase V_t can be seen above
- A rough rule of thumb is that the cell V_t changes by 1.0V/dec under erasing conditions (~100us-1ms).
- Actual NAND Product uses a +ve Well voltage with grounded gate to erase the cell. However at Param (ETEST), often a -ve voltage is applied to the WL and Well grounded. Electrically both are equivalent

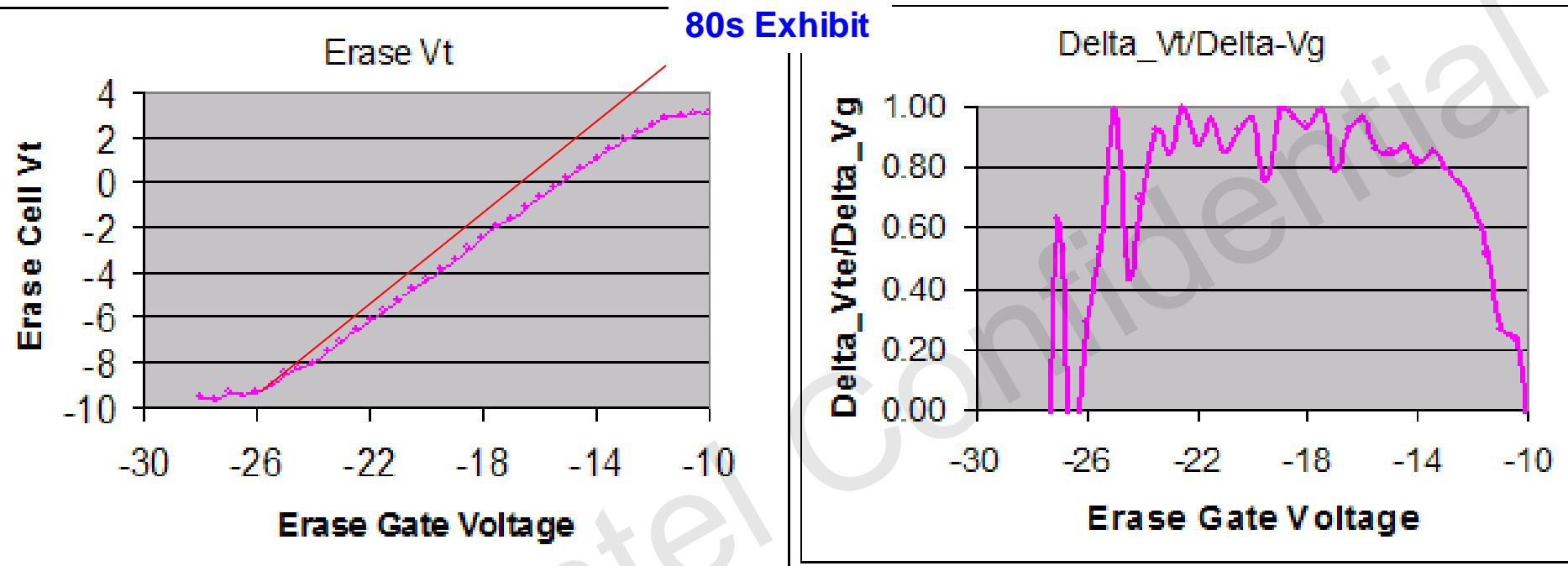
Erase Inhibit (2D NAND)

- We described the NAND Erase as FN tunneling from FG to the channel WL grounded and the P-Well high at 20V. But the P-Well is common to all the blocks
- So, the question is how do we prevent (“inhibit”) all the other blocks from erasing as well, while erasing only the selected block.
- The way we achieve inhibiting Erasing on all the other blocks is called the “Erase Inhibit”

- P-Well goes high for all blocks
- WL on Selected blocks grounded, so these get erased
- Erase inhibited in deselected blocks by floating the WL & allowing them to get “boosted” up along with the Well



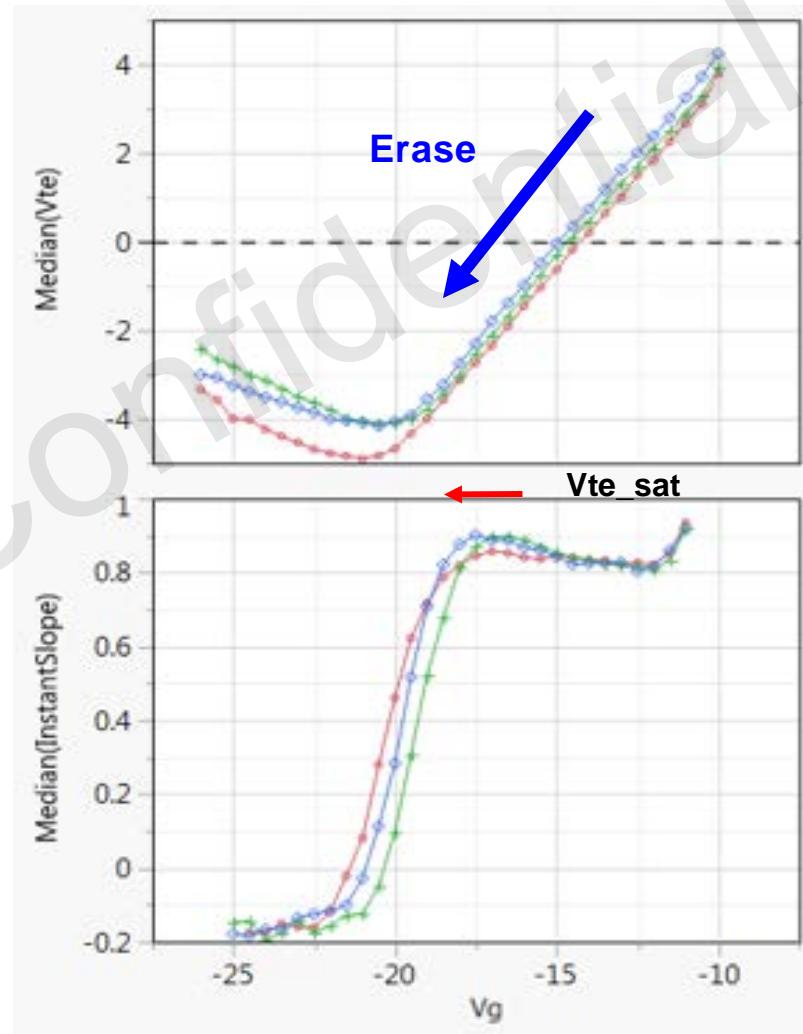
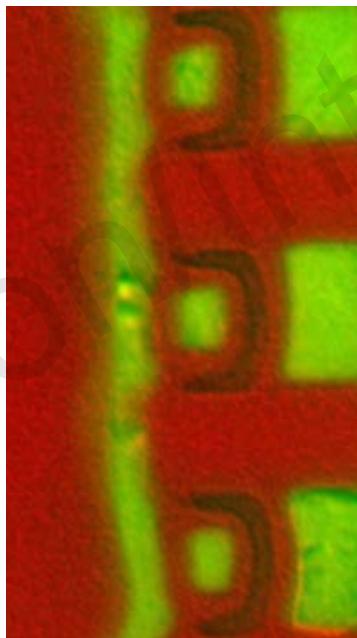
Erase Vt vs Vg (Vw)



- For a given Erase time, the Erase Vt tracks the erase voltage. V_wV_t ($= V_w + V_{te}$) or V_gV_t ($= V_g - V_t$) is a good cell erase metric.
- V_gV_t or V_wV_t is primarily determined by the Cell UV_Vt, FN Tunneling coefficients, Tunnel-oxide and ONO thicknesses, GCR, etc.
- We expect $\Delta V_{te}/\Delta V_g$ to be ~1.0. However, deviation from 1.0 is seen due to GCR difference between Erase and Read conditions as well as Poly depletion effects.
- At very low erase Vt the Vt saturates because of electrons tunneling from the CG through the ONO into the FG balances the electrons tunneling out.

Erase V_t vs V_g (V_w) (100s)

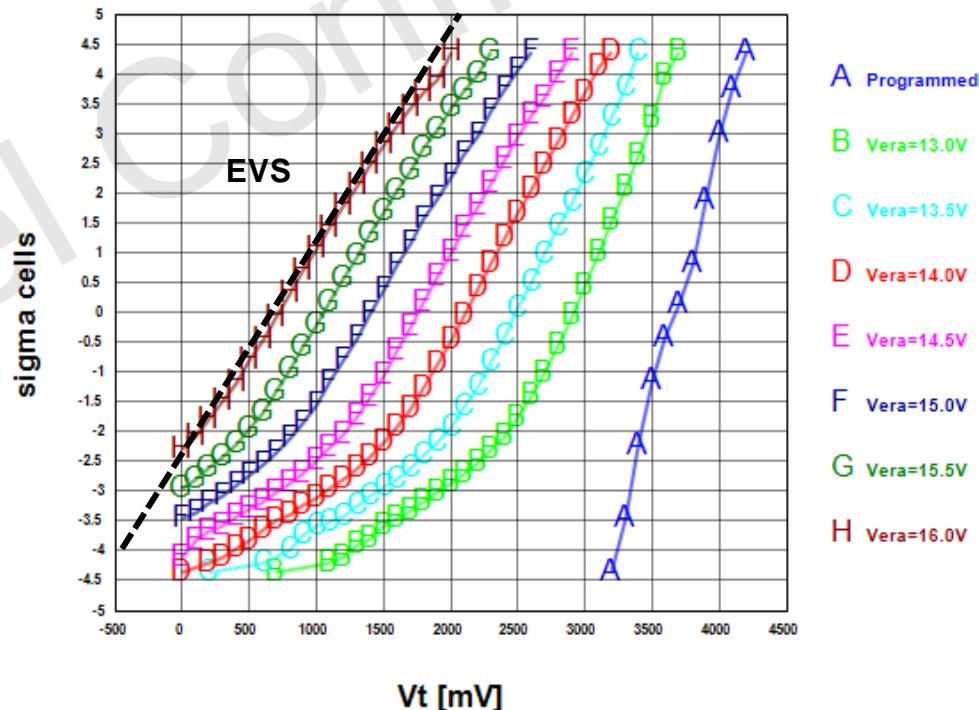
- 100s has shallower erase saturation V_t
- Some of the factors for this are:
 - Lower GCR
 - N+ doped Poly
 - Sharper corner for CG
 - Flank Nitride Shunting the electrical distance to channel



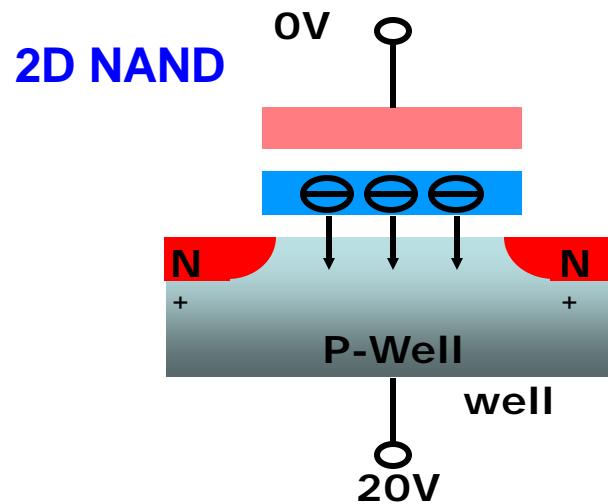
Erase Vt Sigma

Similar to the Program Vt Distribution within the array, we have erase Vt distribution as well, since all the cells in the array do not erase to the same Vt due to cell to cell variations

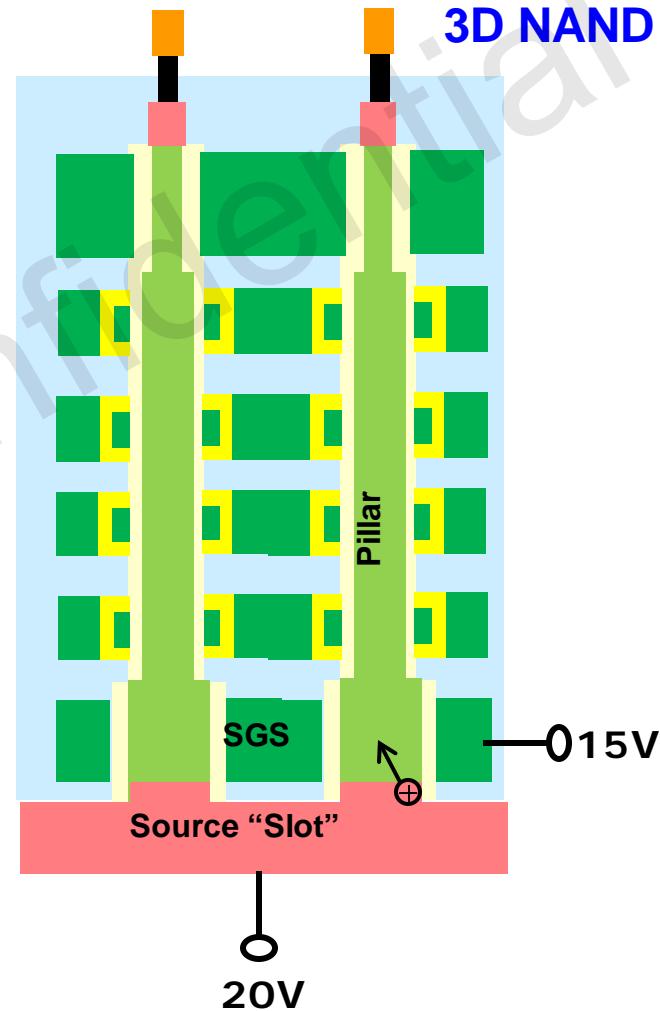
- The erase distribution comes about from differing rates of erase between different cells – (Coupling ratio variations, Z and L variations, UV_Vt variations, FN-Tunneling coefficient variations, Pillar Voltage variation...)
- During erase we need to make sure that the slowest to erase cell is erased well below R1 Vt and it will stay below R1 despite the subsequent FG-FG coupling from neighboring cells



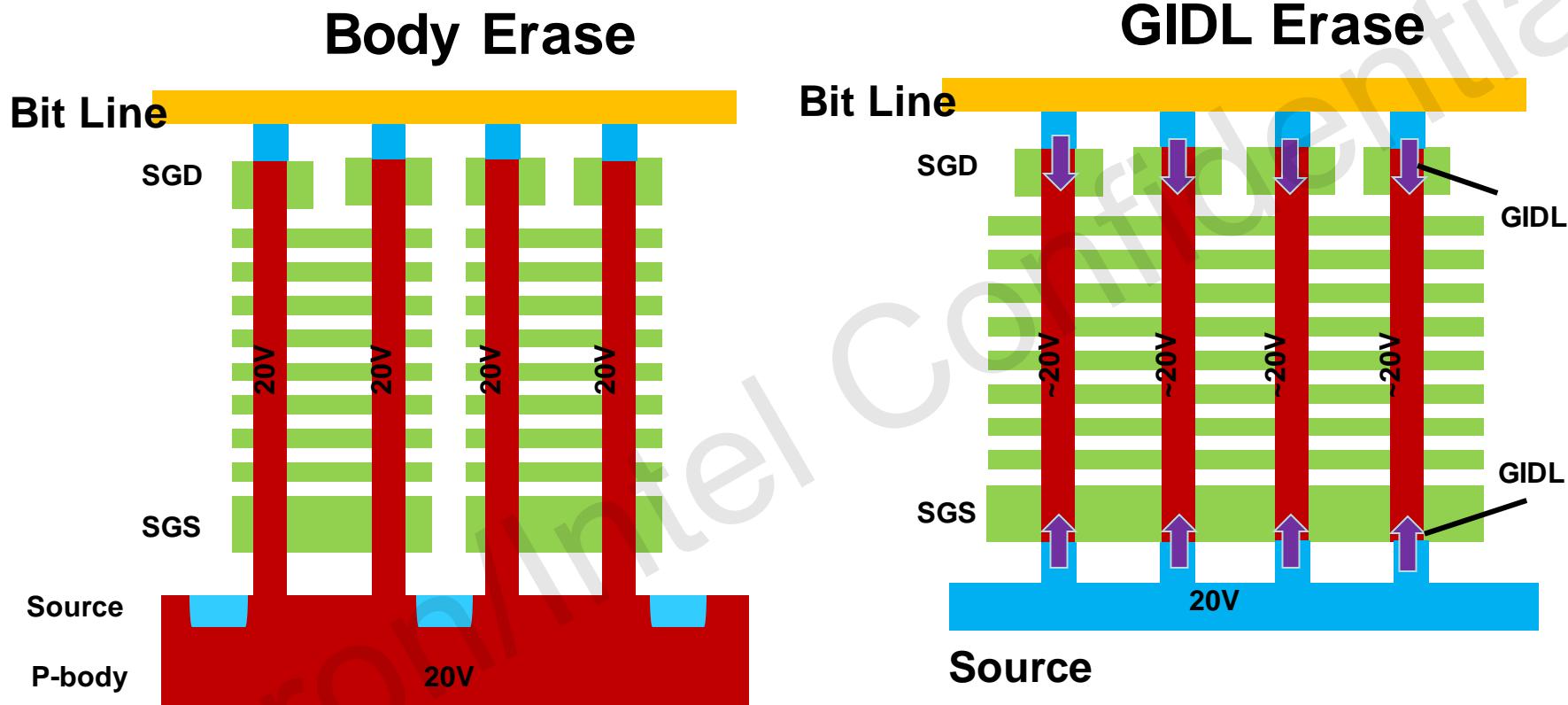
3D NAND Erase



- 2D NAND is erased by applying ~20V to the P-Well/Channel
- 3D NAND requires the voltage being applied to the source slot and transferred to the pillar
- Requires injecting holes into pillar, but n⁺ slot cannot do this.
- Solution: Use SGS GIDL leakage (right)
- SGS reliability from hot-carrier damage is a key concern for 100s 3D NAND.

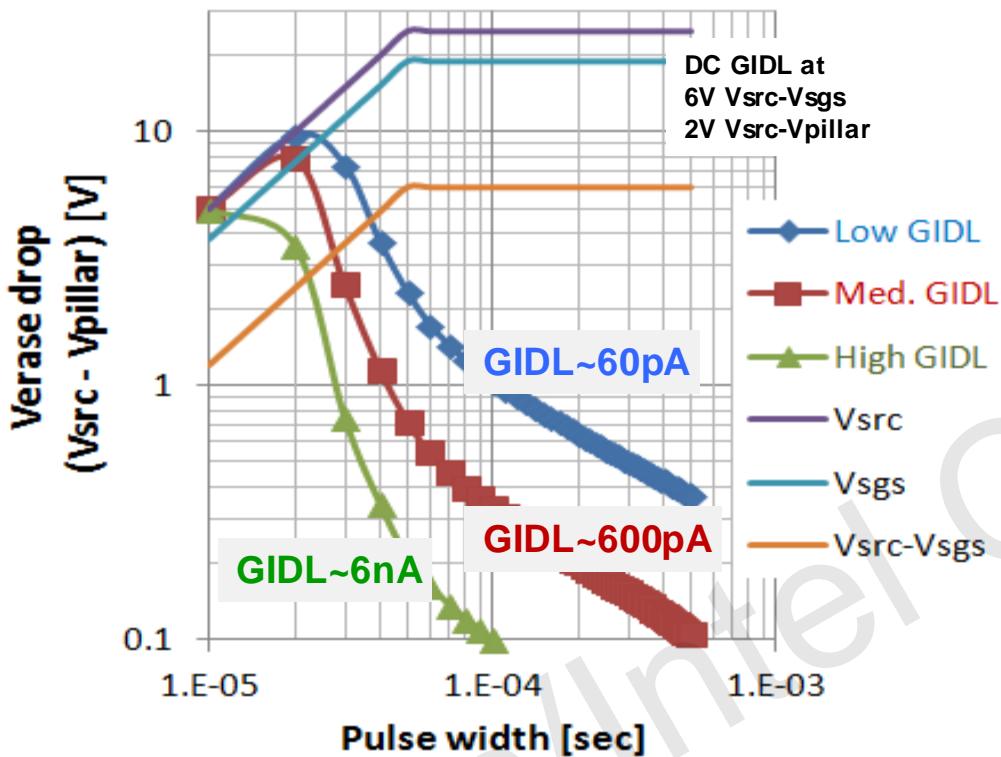


Body erase and GIDL erase

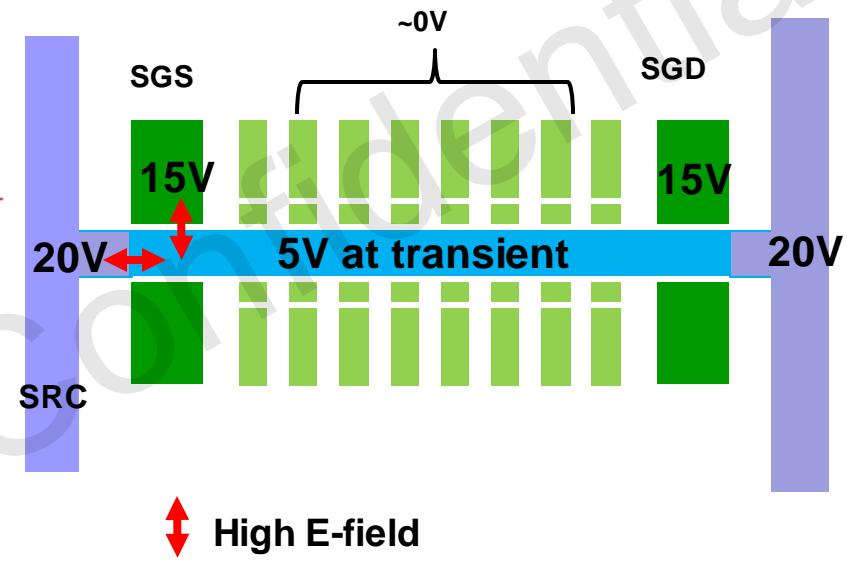


- Body Erase directly biases the pillars at erase voltage.
Pillar to P-body interface resistance control is challenging
- GIDL Erase requires sufficient GIDL current to charge the pillars.
- 100s/110s do the GIDL erase.

Erase GIDL current requirement

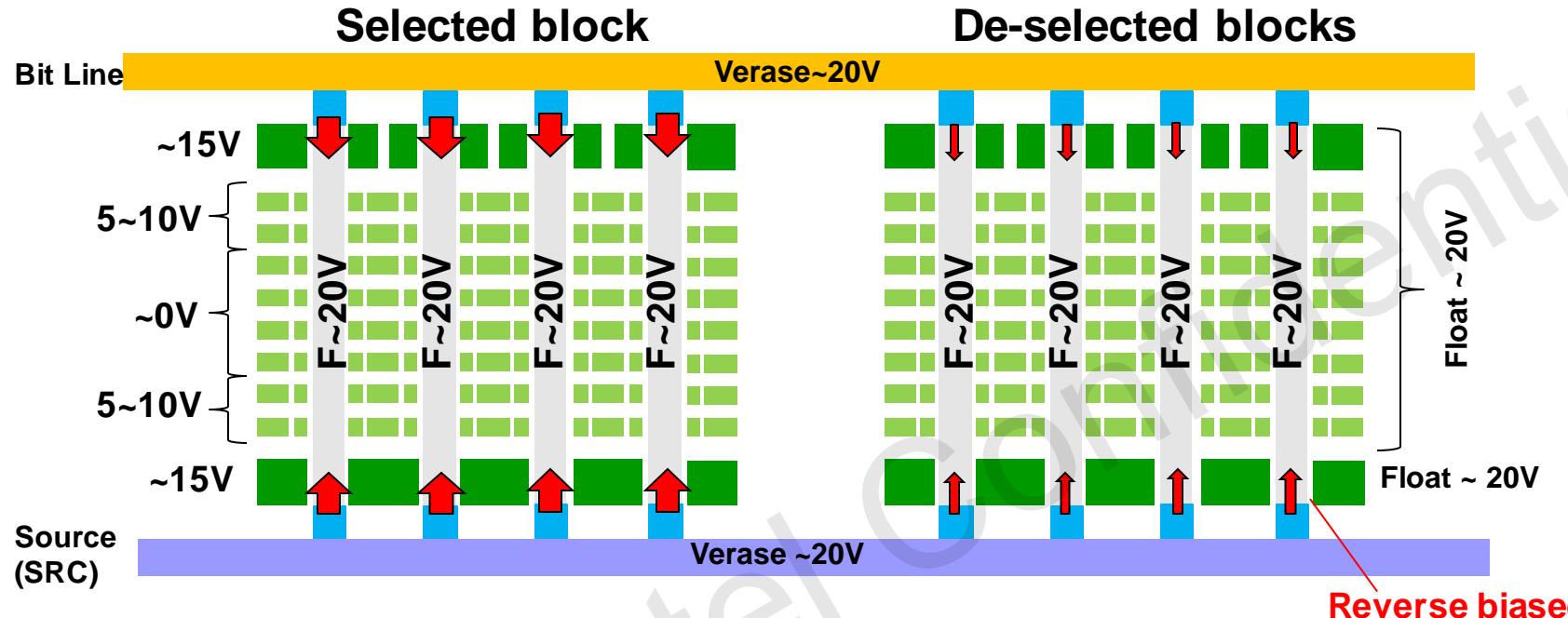


Transient High E_field (and low Vpillar)
due to insufficient GIDL



- SGS GIDL current has to be large enough to charge pillars within the erase ramp time.
- Insufficient GIDL current results in large voltage drop between V_{src} and V_{pillar}. Transient Hot-e and SG ox stress, High Verase and large erase property variation (EVS).

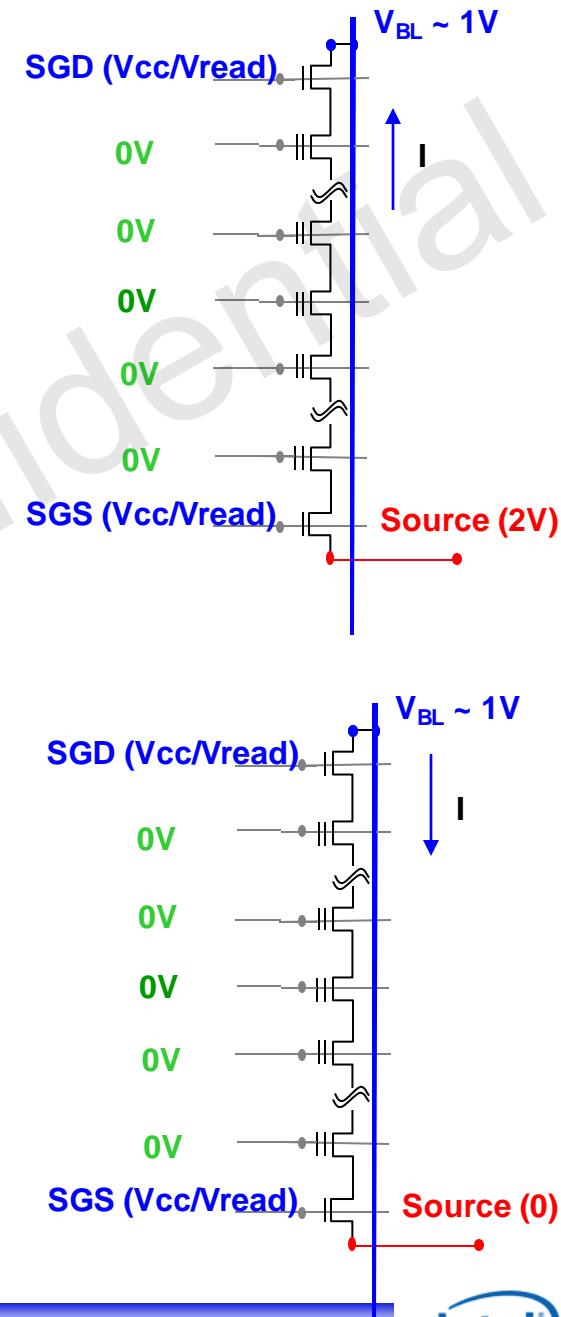
3D NAND Erase Bias Condition



- SRC/BL is biased at erase voltage of ~20-25V
- SGS/SGD is biased ~5V below SRC (negative Vgs) on Selected Block to generate GIDL current to charge up the pillar.
- WLs of the Selected block are grounded to erase them
- SGS/SGD on Deselected block is floated
- WLs of the Deselected Blocks are floated which makes them go up to the same potential as the Deselected block Pillars → Inhibits erase

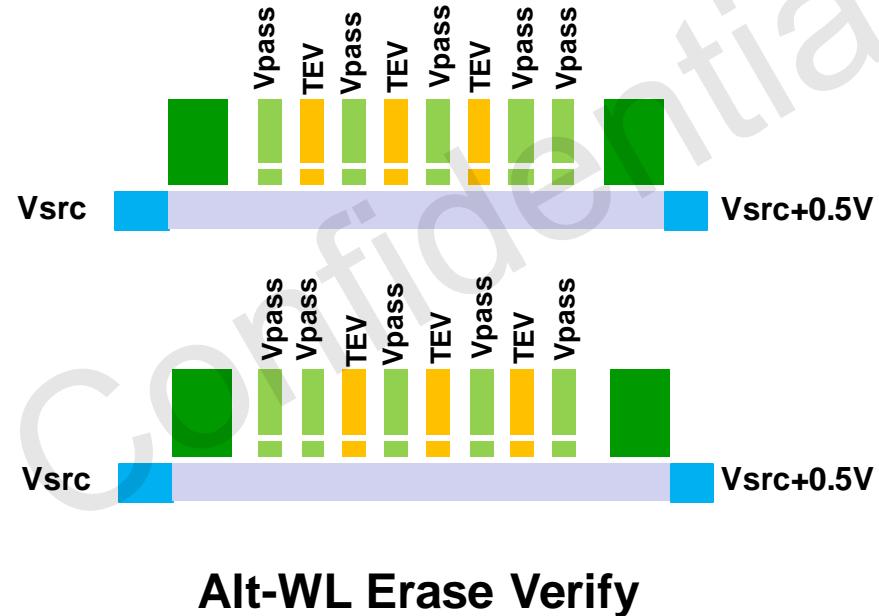
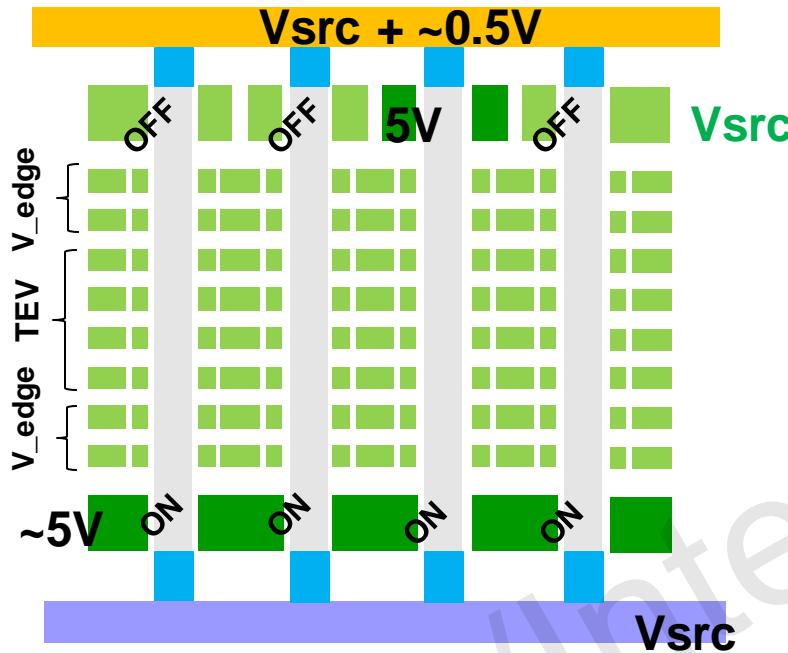
Erase Verify (EV) (2D)

- Historically (until 80s) Source side verify for erase verify (Assure cell $V_t << 0V$)
 - The source is biased to $\sim 2V$. WLs at $\sim 0V$.
 - The bitlines are floating. Due to the current from the source, the bitlines float up.
 - To pass the verify, the V_{BL} needs to go $>1V$. For the BL voltage to float up above 1V, all the cells in the string need to have a " V_t " $< -1V$.
 - Adjacent WL bias is different between EV (~ 0) and Read ($V_{passR} \sim 7V$). This allows EV to verify negative V_t by the amount of $(V_{passR_n} \pm 1 \times WL-FG\ coupling \times 2sides)$.
 - The body bias on the cell provides EV to verify negative V_t by additional amount
- Conventional Drain side verify can be used as well and rely on WL-FG coupling to achieve Read $V_t \ll 0V$.



3D NAND Erase Verify

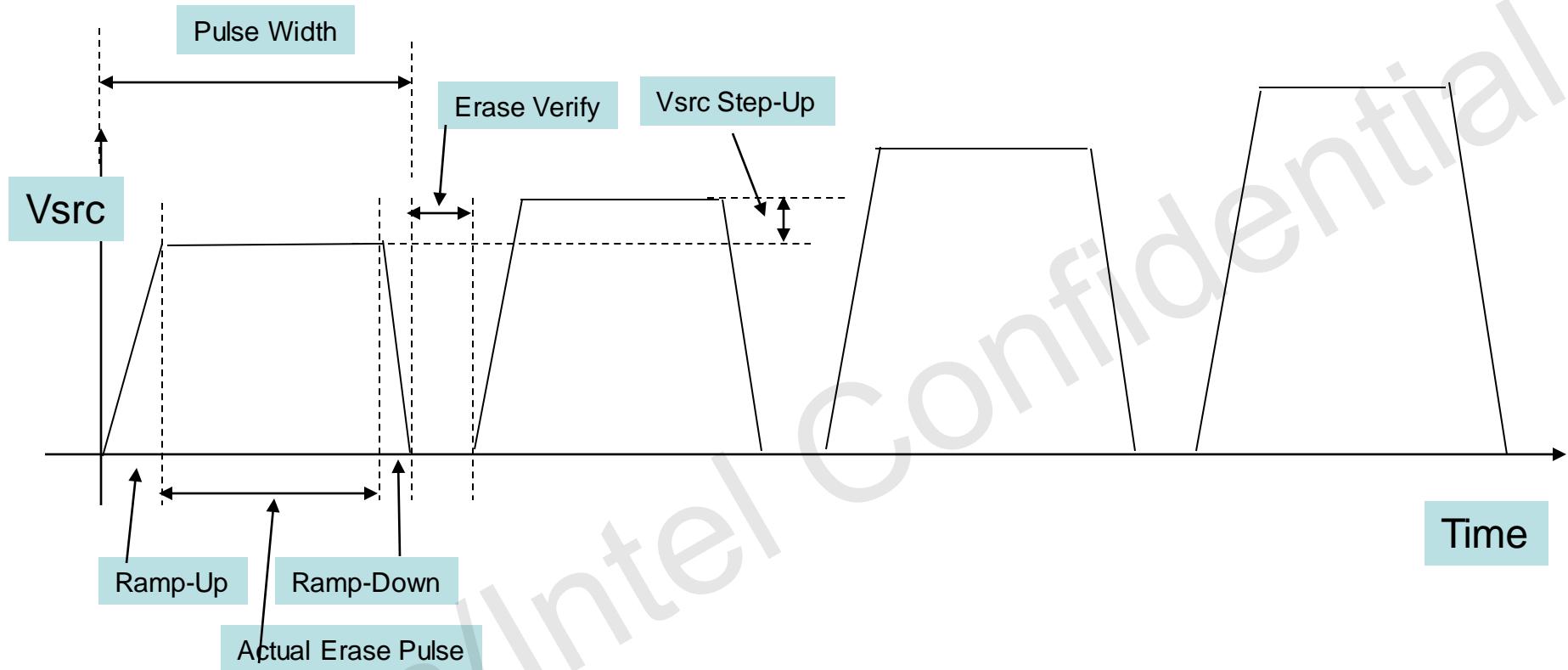
Erase Verify



Alt-WL Erase Verify

- Each sub-block needs to be individually verified.
- Lack of LDD in the 3D makes the string shut-off at low V_g (no inversion charge over the LDD at low gate bias) → Simple string verify with TEV on all WLs do not work. → Two pass verify for each string with every other WLs at V_{pass} (~8-9V) to invert the LDD region.
- Continue with erase pulses until all sub-blocks pass verify ($Cell\ V_t < Tgt$)

Erase Voltage Pulse Shape



- Since Erase V_t is once again a logarithmic function of time, 1V of additional erase will require 10X longer erase time or 1V higher Erase voltage.
- By using an erase voltage that is stepped up after each verify, the exponential time impact can be turned into a linear time impact

Short Read Window Budget (RWB) Overview

Full (nearly) RWB Component List

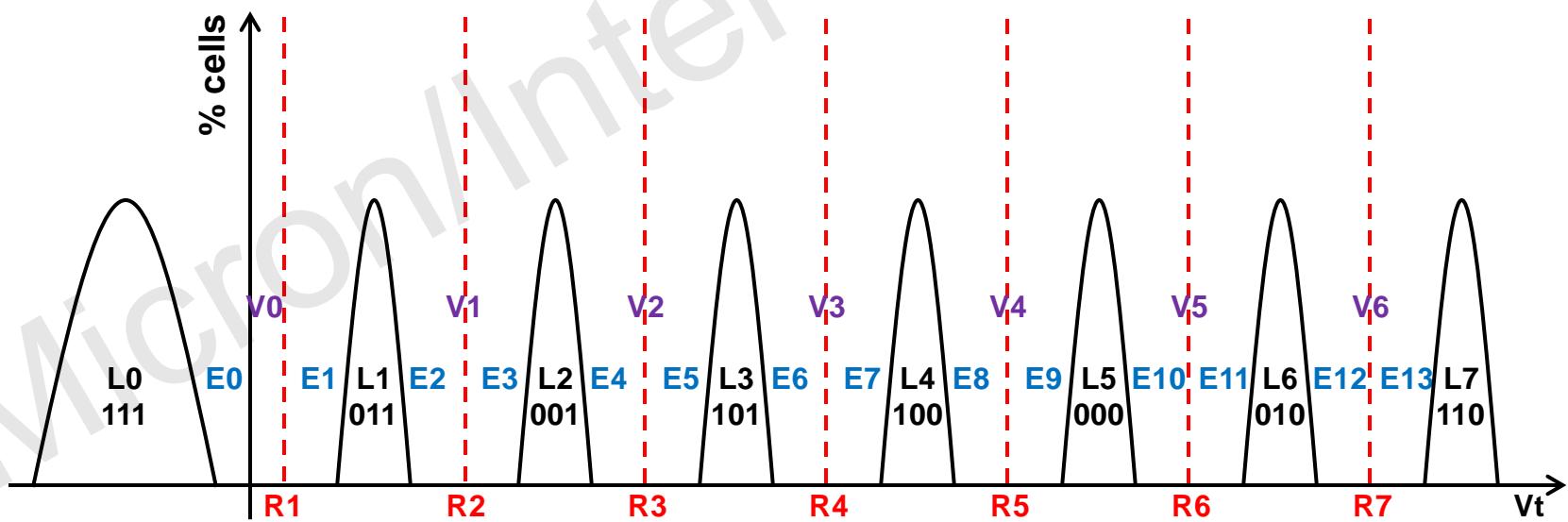
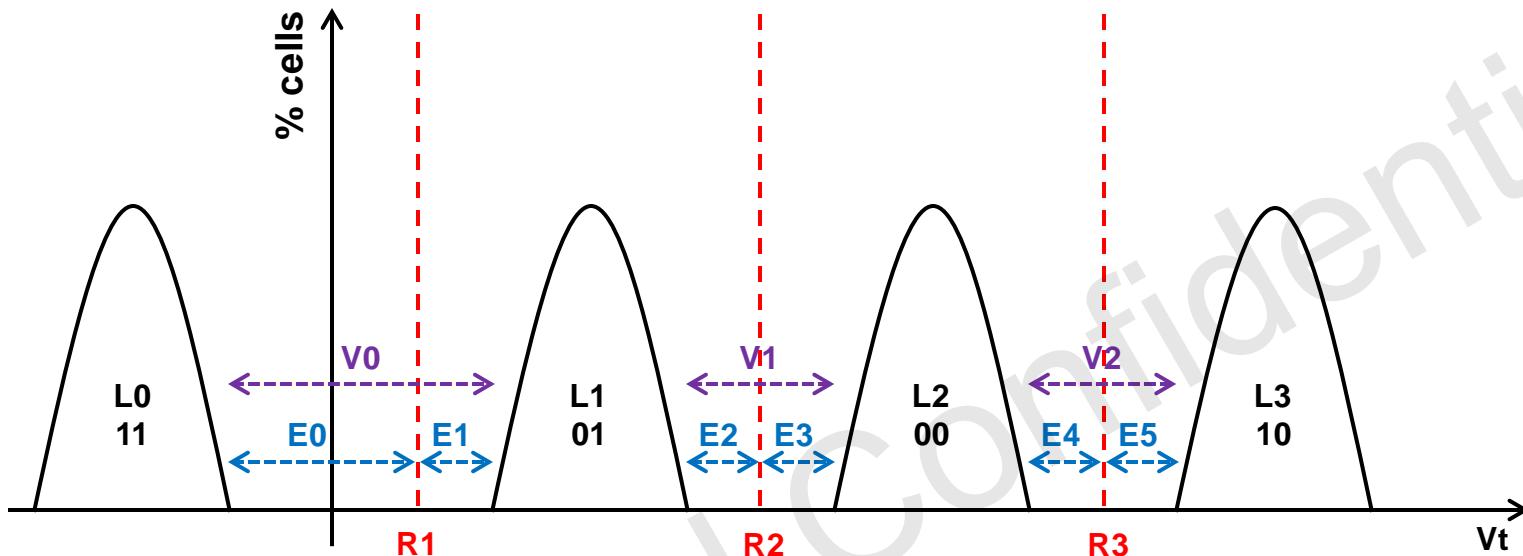
Cell/TD/PE Parameters:

- RTS
- DRB/ICL
- FG Coupling
- VgVt Sigma
- Erratic Programming
- Back-Pattern
- Source Rise in Verify
- Poisson Pgm Noise
- Read Disturb
- Cell Tempco Sigma
- Pkg Stress
- Process Variation
- Trim Errors at Probe
- Accounting for all these factors to arrive at the bottoms up the different state-widths and the margin requirement is the RWB methodology
- Post silicon we can directly measure the silicon and confirm the net window

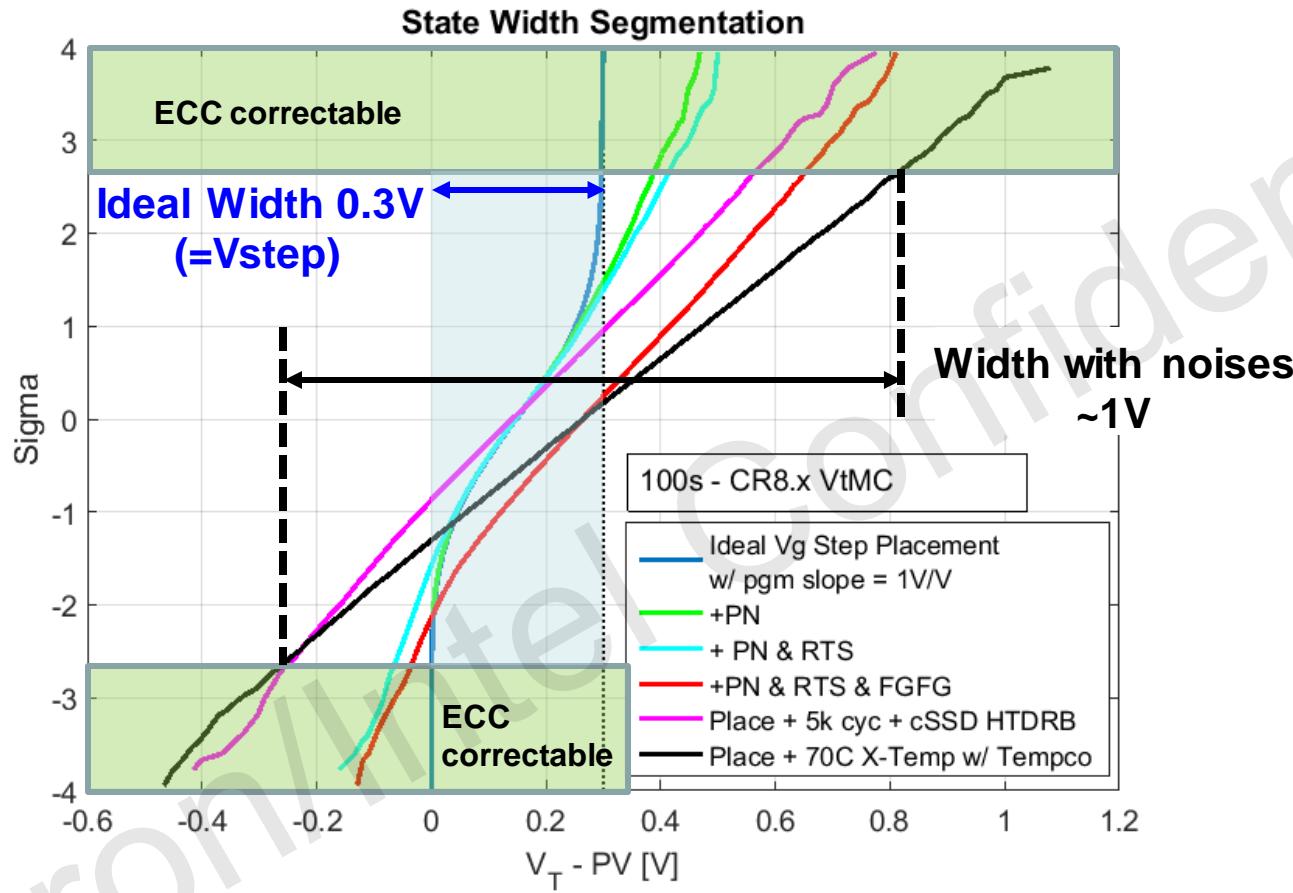
Alg/DE/Spec Parameters:

- ECC Capability
- Pgm Gate Step & SSPC
- Sensing Noise
- Data Pattern (PRP, etc)
- DAC Range & Step
- Temp Compensation Sigma
- Sense Amp Offset Sigma
- Circuit Response to Pkg Stress
- WL Voltage Regulation
- Random Unit Variation
- Density & Architecture
- Cycle Count
- Read Count

MLC/TLC distributions



State Creation 100s Example

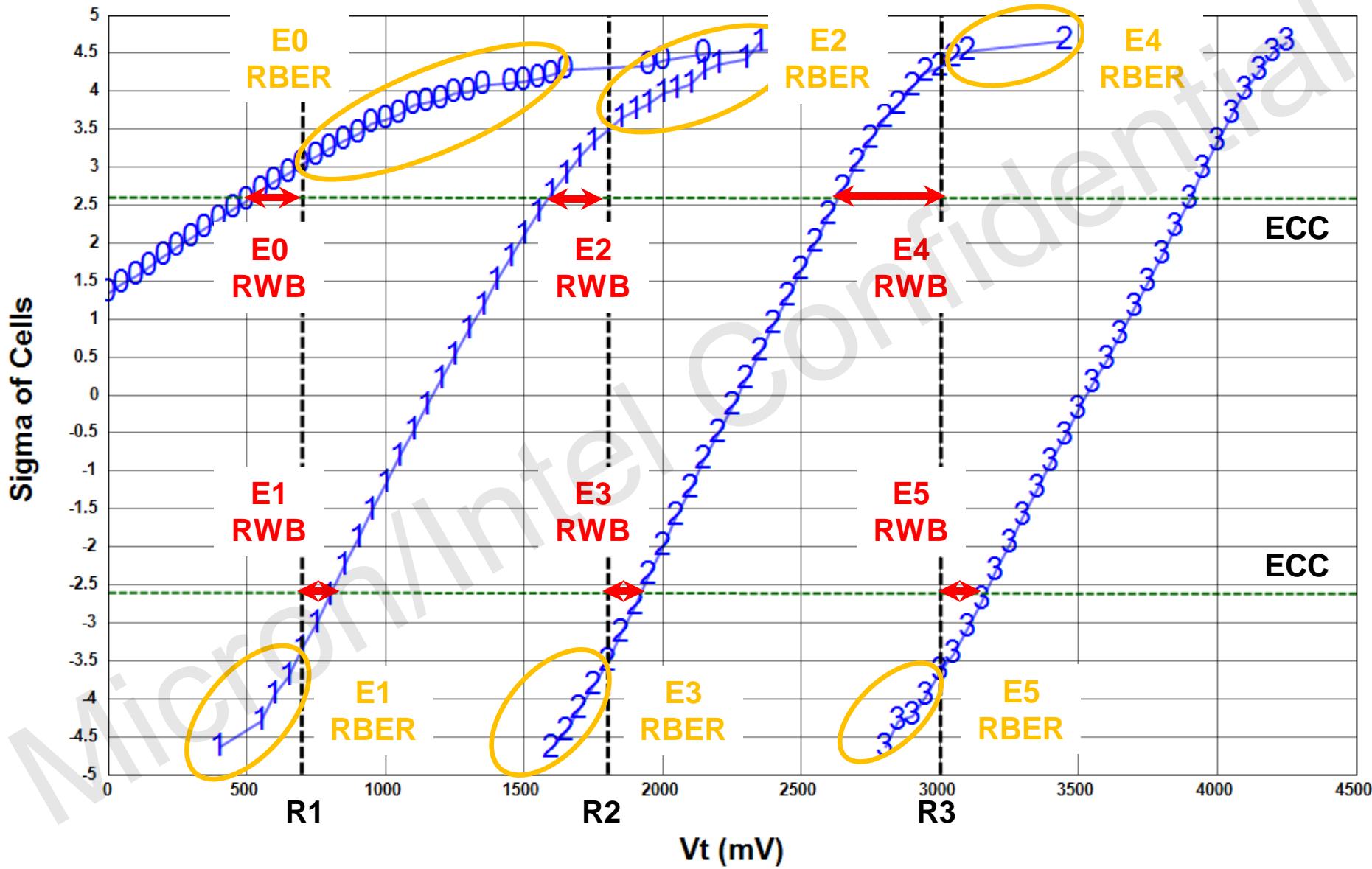


Assumptions: 300mV gate step, coupling from 1WL neighbor, MLC 1-pass algo.

Actual Vt width with noise components

is much larger than the ideal width dictated by Vgstep only.

RWB and RBER definition

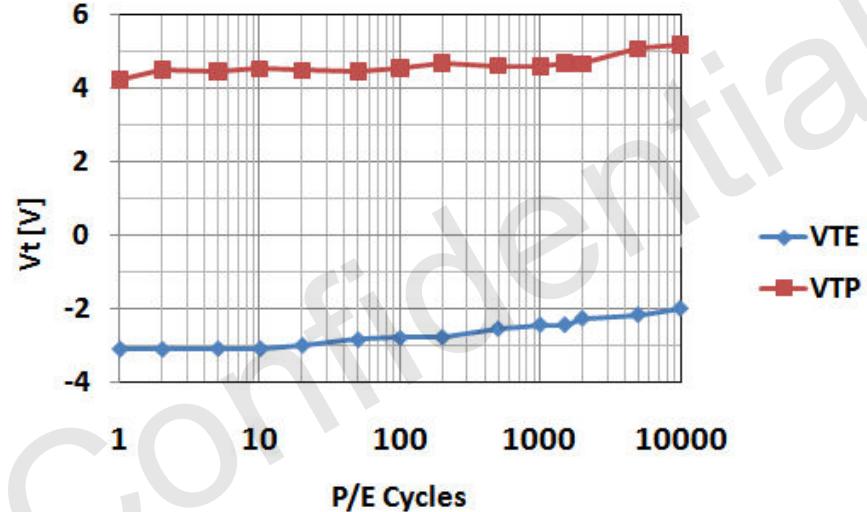
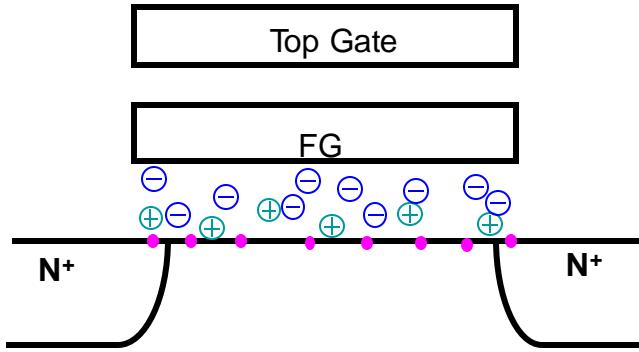


Short Cell Reliability Overview

Reliability Overview

- Endurance & Trap-up
- Data Retention – SBCL, ICL
- RTS
- Read Disturb – FN & Hot-e
- Transient Vt and First Page Read effects
- Floating Body Effect
- ECC

Endurance



- Tunnel Program/Erase generates uniform **electron traps**, **hole traps**, and **interface states**
- Trapping of electrons in the Tunnel oxide lead to an increase in the threshold voltage of the Cell. Typically Program V_t and Erase V_t increases (at fixed P/E voltage/time). Erase V_t shift is typically $>$ Program V_t shift.
- Need to make sure that even after cycling we are completing erasing all the cells to V_t well below R1 level.

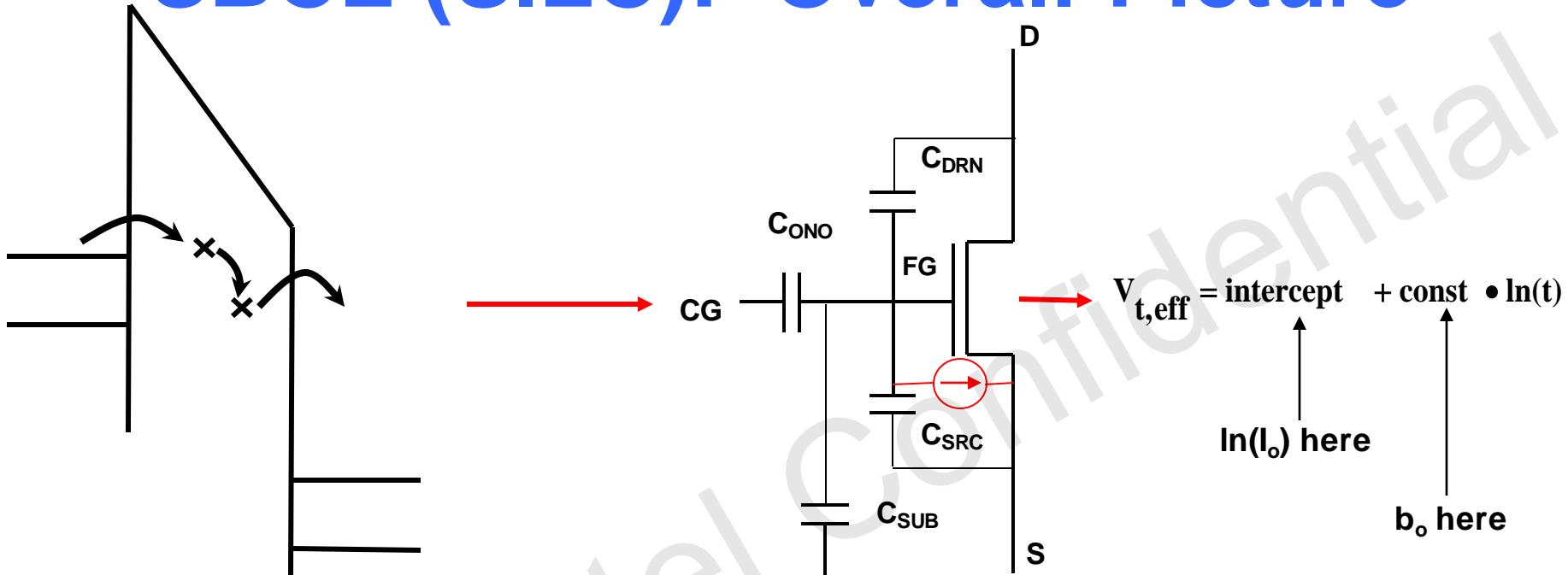
Trap-up impact on reliability

- Lower VgVt increases the sensitivity to the pass voltage during read, degrading read disturb
- Increased PVS (i.e. wider distribution) degrades PD because inhibited distribution gets wider
- Very fast to program cell can get misplaced after the first program pulse or due to hot-electron injection (over-programming)
- Channel interface degradation can trigger lower capability to hold boosting voltage further degrading PD (especially at high temperature)
- Trapped charge contributes to the cell Vt during program operation, but is released more easily over time, impacting RWB odd edges (E1/E3/E5) due to Intrinsic Charge Loss (ICL)
- Higher VwVt leads to lower endurance capability
- Also increased Tev-Vttop distance leads to lower endurance capability

Data Retention

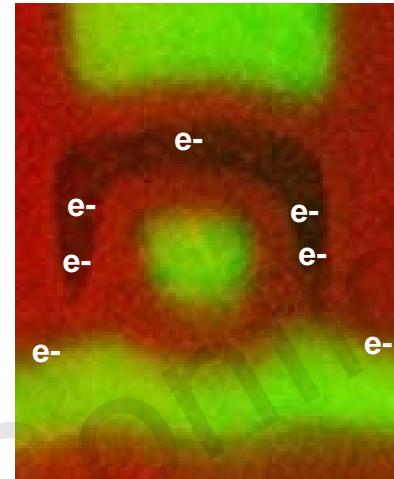
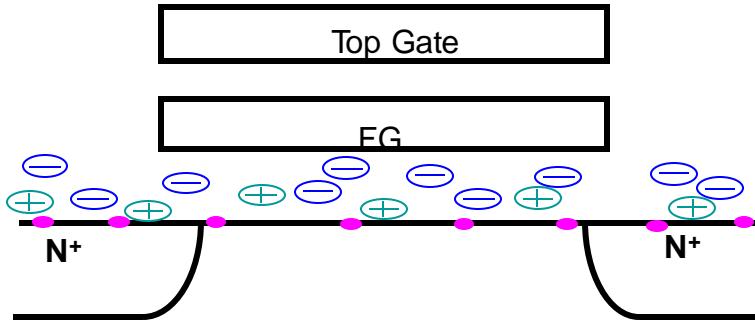
- Common Data Retention issues include:
 - Single bit charge loss or charge gain due to SILC related mechanisms
 - Cell V_t shift due to the de-trapping of trapped charges in the oxide
 - Charge loss or charge gain due to oxide defects, physical defects, contaminations, etc.

SBCL (SILC): Overall Picture



- Cycling leads to traps. Traps lead to TAT
- TAT has an exponential I-V, with a magnitude exponentially dependent on trap position
- Random trap positions lead to a wide defect tail which moves in parallel fashion with $\ln(t)$
- Individual-cell SBCL is noisy because broken bonds are erratic – this is known from many oxide mechanisms
- The variation that we care about is in the intercept term, which stems from the I-V prefactor
 - Cell-to-cell variation is extreme value (in V_t), Weibull in time => based on our data and on STM's Monte-Carlo sims and consistent with percolation simulations of oxide TDDB
 - Cycling has ~ square root dependence on defect density, as do many other oxide mechanisms
 - Temperature dependence is weakly positive, consistent with other tunneling measurements in SiO_2
 - Defect density increases with hole injection, consistent with other oxide degradation mechanisms
 - Thicker oxide has exponential effect, ~consistent with 10x per 2 Angstroms in WKB ($10(1/6) = 1.5x/\text{Angstrom}$ predicted)
 - Changes in UV V_t are predictable based on the equivalent circuit – simple V_t offset in asymptotic region

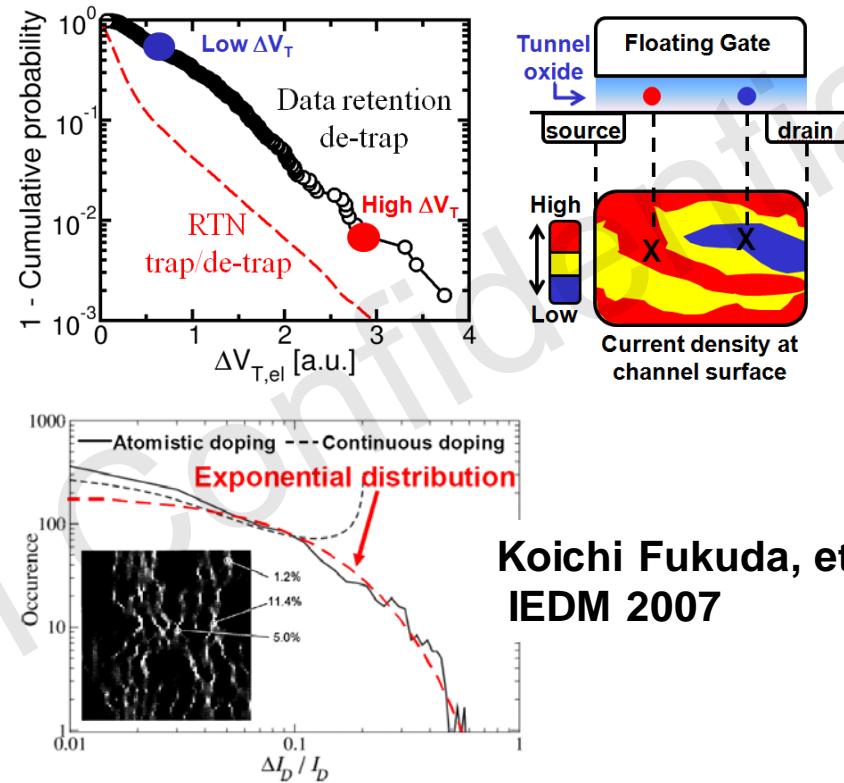
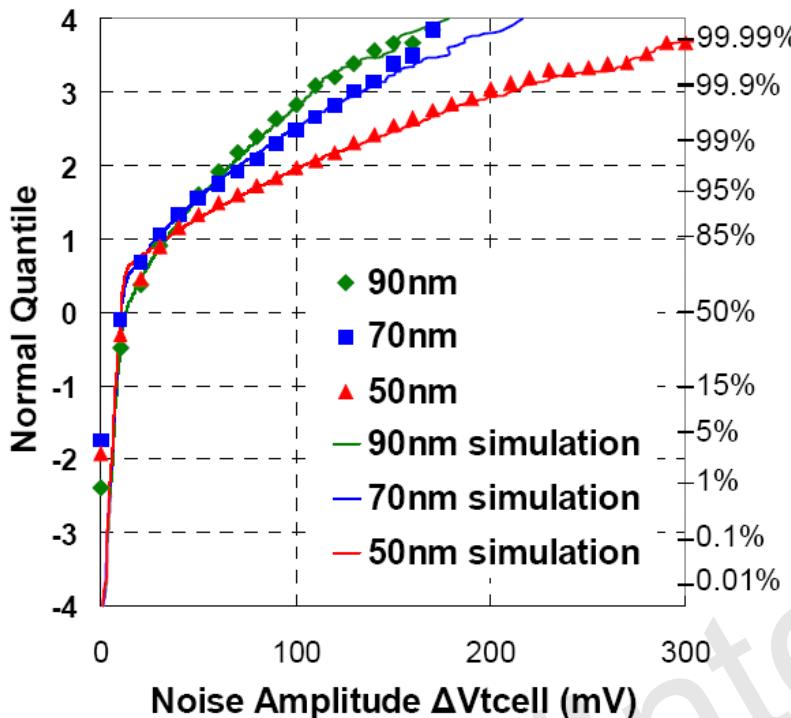
Intrinsic Charge Loss (ICL)



Parasitic trap/de-trapping Location in 3D NAND.

- Tunnel Program/Erase generates **electron traps**, **hole traps**, and **interface states**
- The trap-up occurs in the channel as well as the S/D as well as nearby dielectrics and all these trapped charges can effect the cell V_t. Detrapping of these will cause the V_t to shift
- Most Window closure and post cycling high temperature bake charge-loss (ICL) are due to bulk electron de-trapping
- ICL is also called HT-DRB due to the high activation energy causing the problem to be severe at high Temp

Random Telegraph Signal (RTS)

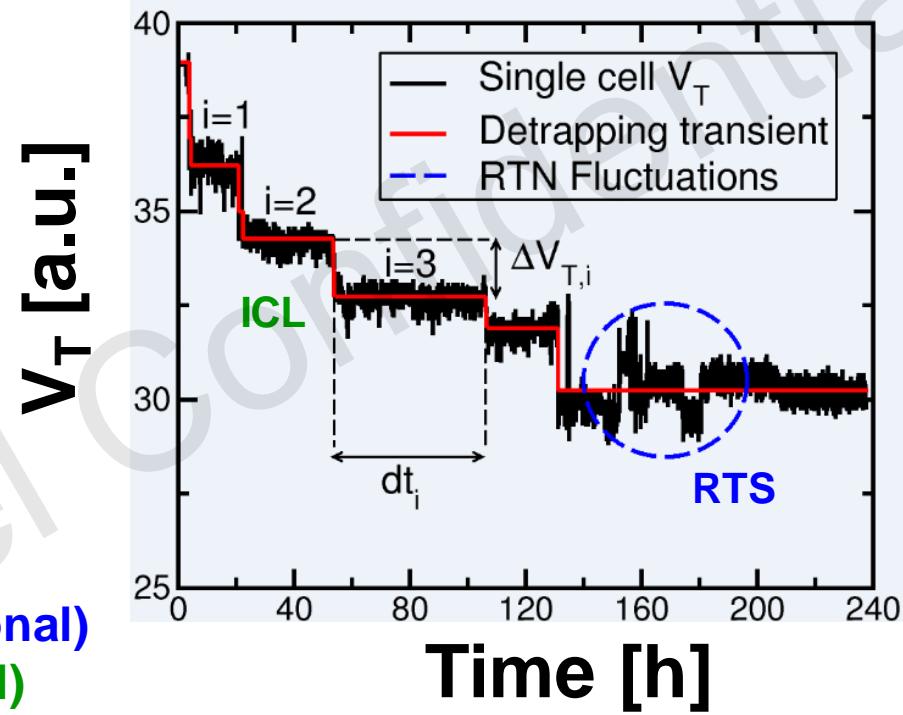
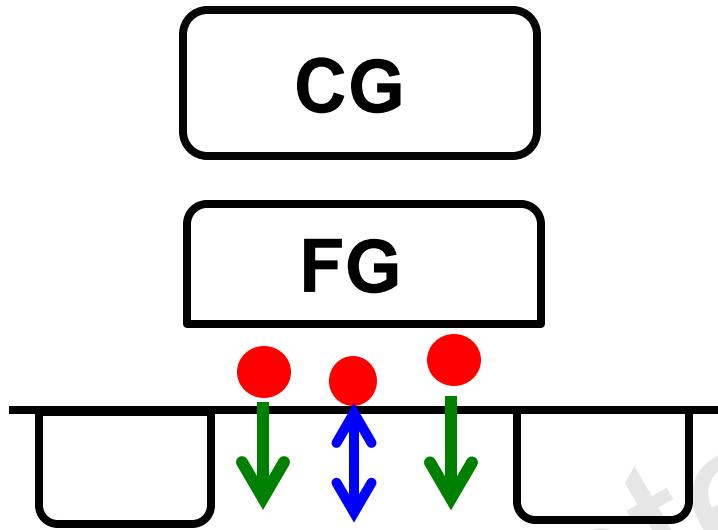


Koichi Fukuda, et al.
IEDM 2007

- Large-amplitude-noise tail portion of the noise distribution produced by the current percolation effect
- The noise amplitude is found to follow $1/(L_{eff}W_{eff})^n$ making it worse as the cell is scaled
- In 3D w/ Poly channel, the additional grain-boundaries in the current path will impact RTS. However, large W of 3D reduces the impact.

ICL and RTS

Time dependent V_T fluctuations

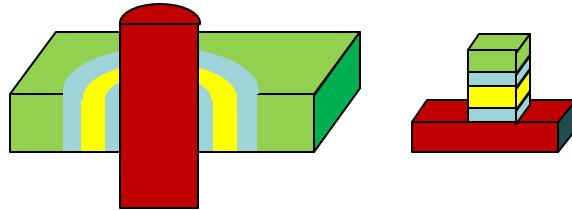


RTS (Capture/Emission, bi-directional)
ICL (De-trap, mainly uni-directional)

Both of ICL and RTS are V_T shift caused by charge trapping/de-trapping mainly at tunnel oxide. V_T shift characteristics and statistics can be described by

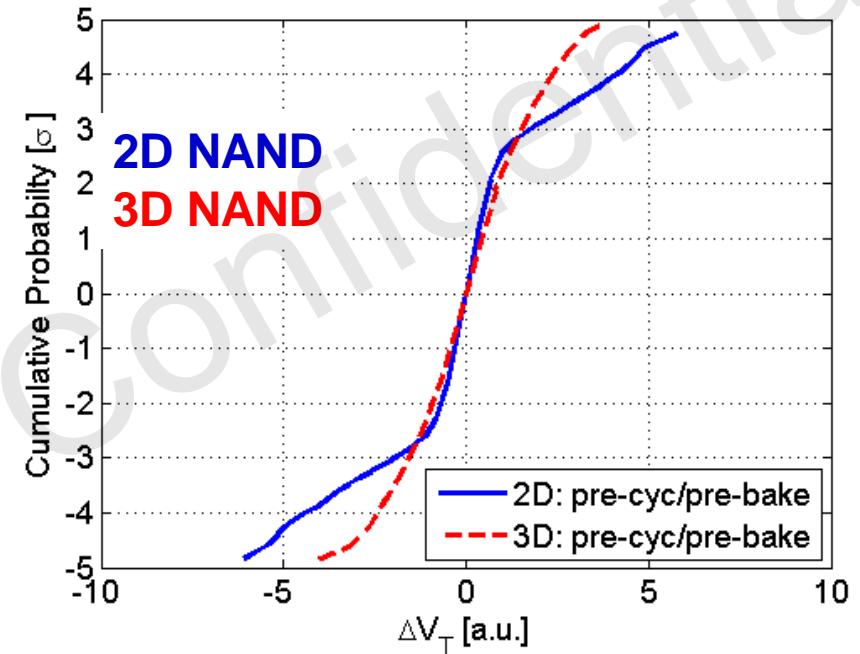
- (1) # of traps and trap/de-trap time constant
- (2) ΔV_T shift per trap

3D NAND vs. 2D NAND



	3D	2D
$\Delta V_T \text{ ave} / e$	Small	Large
#e fluctuation	Small %	Large %
$\Delta V_T / e$ spread by percolative conduction	Large chan. width	Narrow width with STI edge
	Poly Si chan.	Dopant fluctuations

RTS dVT cell by cell histo



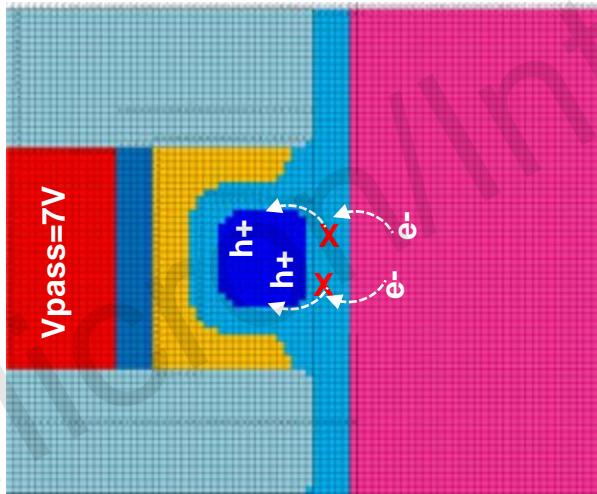
Key difference from 2D NAND is ΔV_T / electron statistics.

With a successful poly Si engineering, better RTS than 2D NAND, owing to the large channel width minimizing percolative current conduction,

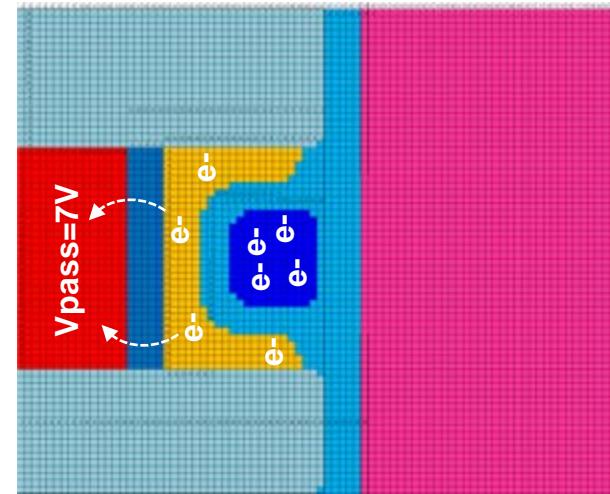
Read Disturb

- Unselected WLs in the string are biased at Vpass ~7-8V on 100s
- This leads to two issues:
 - Charge can be injected through the tunnel oxide, with a TAT mechanism that is the opposite of SBCL, mainly affecting L0 cells, due to higher electric field on the tunnel oxide: RD CG
 - Charge trapped in the IPD2 layer can be ejected through IPD3 to the CG, mainly affecting highest level (L3 in MLC, L7 in TLC), due to the highest charge trapped there and the highest electric field through IPD3: RD CL

RD CG: Read Disturb Charge Gain



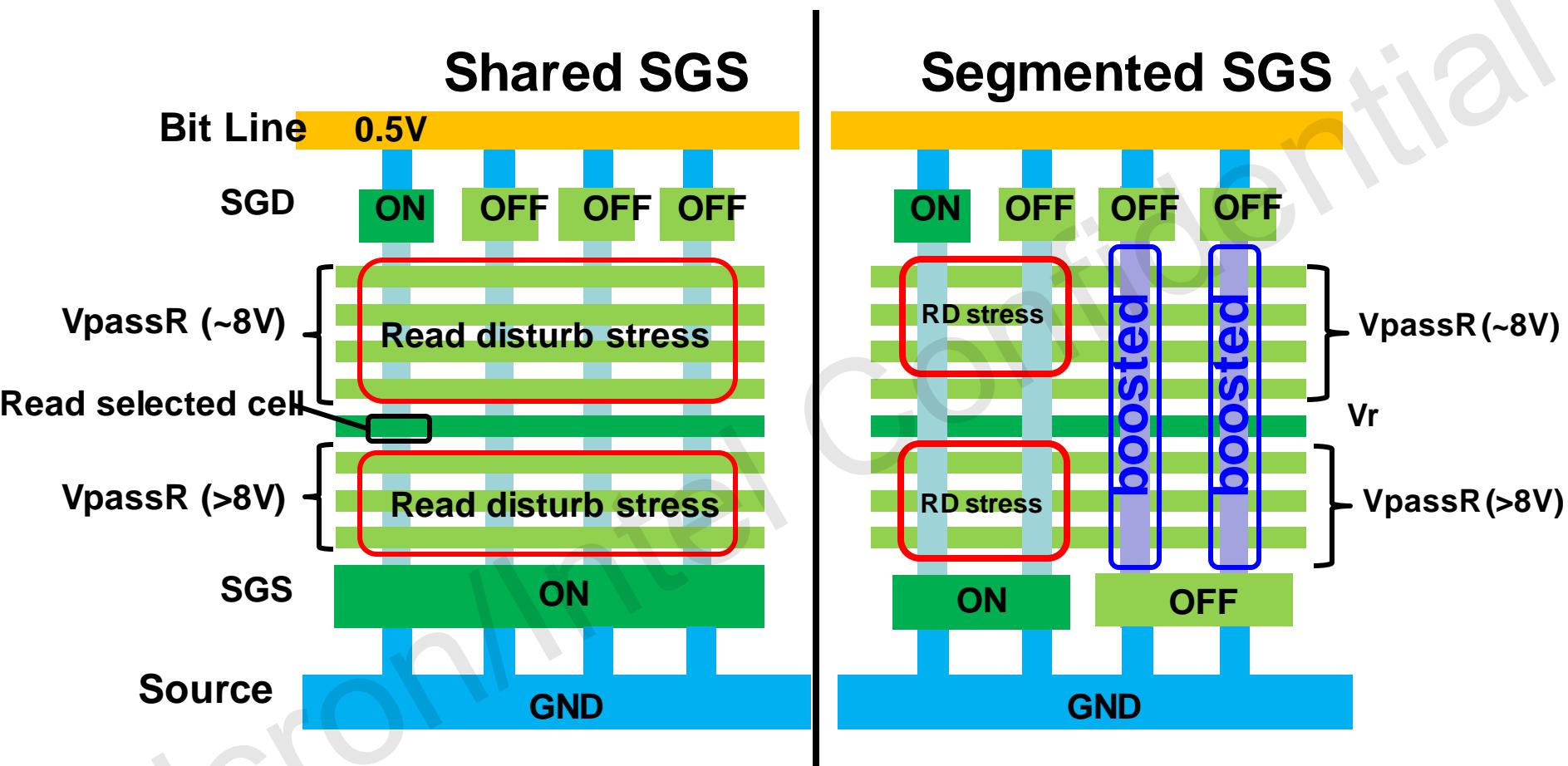
RD CL: Read Disturb Charge Loss



Read Disturb

- Read disturb CG is one of the highest risk mechanisms on 100s
- Vpassr is higher than old technologies (~7-8V vs 5-6V) due to polysilicon channel and missing LDD region: need to create the channel not only under the cell, but also create inversion in the region between two cells
- Also the higher number of pages in the block is impacting, since the reliability goal is defined in terms of number of full block reads
- Read disturb CG gets worse with cycling, due to more traps present in the tunnel oxide (same as SBCL), and also due to higher UVVT of the cell (opposite of SBCL)
- Higher Vpassr makes also FN component present on 100s, together with TAT, creating a dependence on the VgVt of the cell
- Goal is defined at the highest cycle count for the product, e.g. 3k full block reads after 1.5k P/E cycles for TLC, or 10k full blocks reads after 3k P/E cycles for MLC

3D NAND Read Disturb

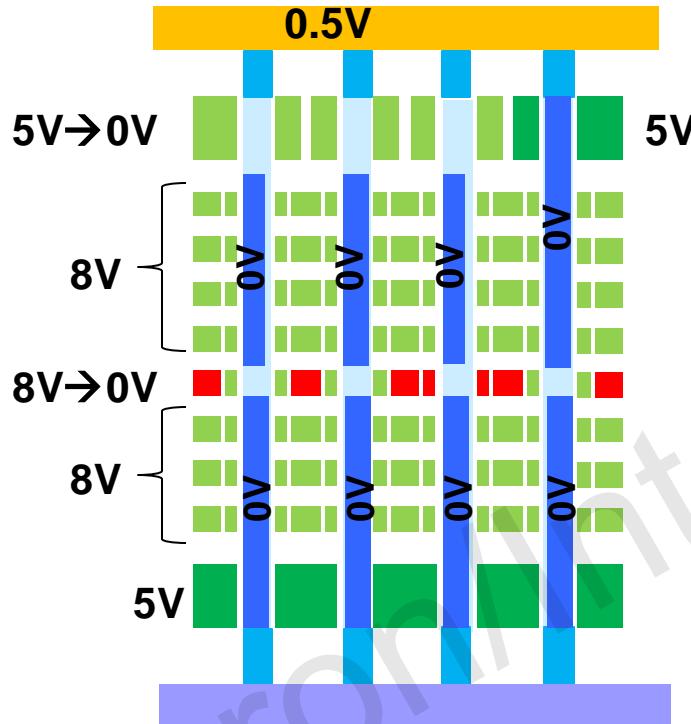


[Shared SGS] Pillars of the de-sel. sub-blocks are grounded, leading extra read disturb stress (Rel. deg.) and WL-Pillar capacitance (Perf. and power loss).

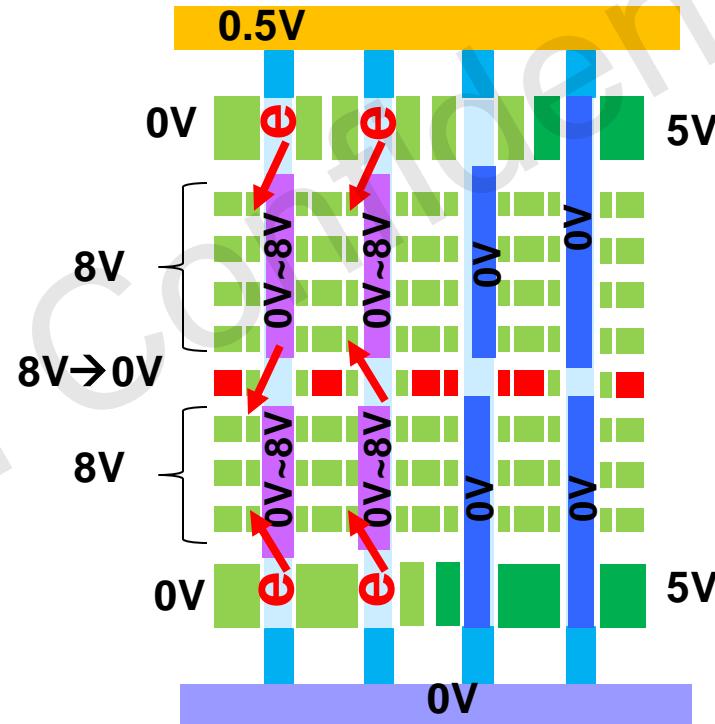
[Segmented SGS] Some of the de-sel. pillars (75% for B1MA) are boosted During read, reducing read disturb stress and WL-Pillar capacitance.

Hot-e RD

Shared SGS



Segmented SGS

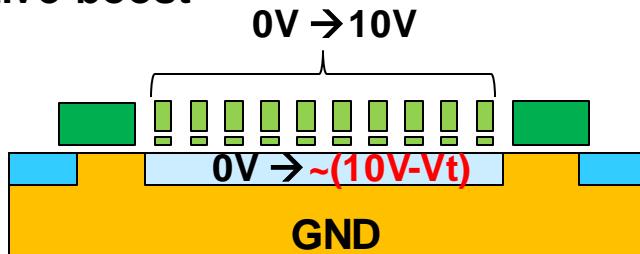


Due to the High E_field at the edge of the boosted pillar,
Hot-e injection disturb can occur in the segmented SGS,
requiring bias and waveform optimization to minimize the hot-e disturb.

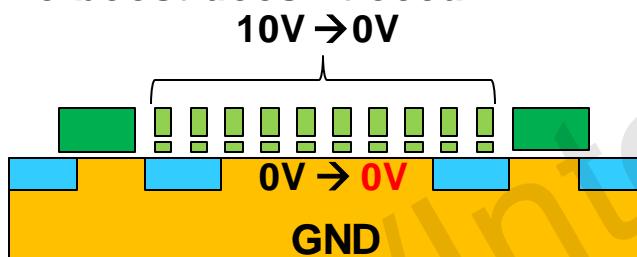
Floating body effect

2D NAND

Positive boost

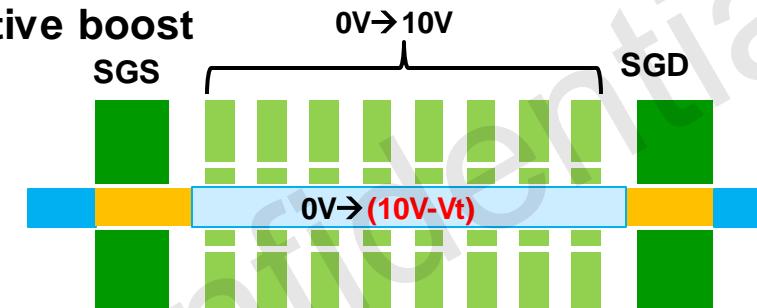


Negative boost doesn't occur.

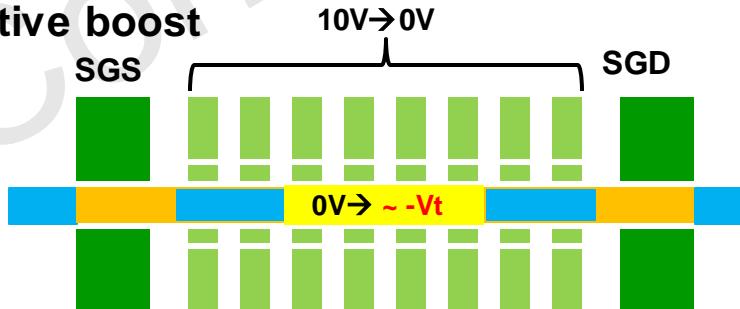


3D NAND

Positive boost



Negative boost



Assume WL to Ch coupling ratio ~1

Lack of electrons vs.
thermal equilibrium

Lack of holes vs.
thermal equilibrium

- **Positive boost can occur in both 2D and 3D NAND.**

This is known as boosting. Lack of electrons.

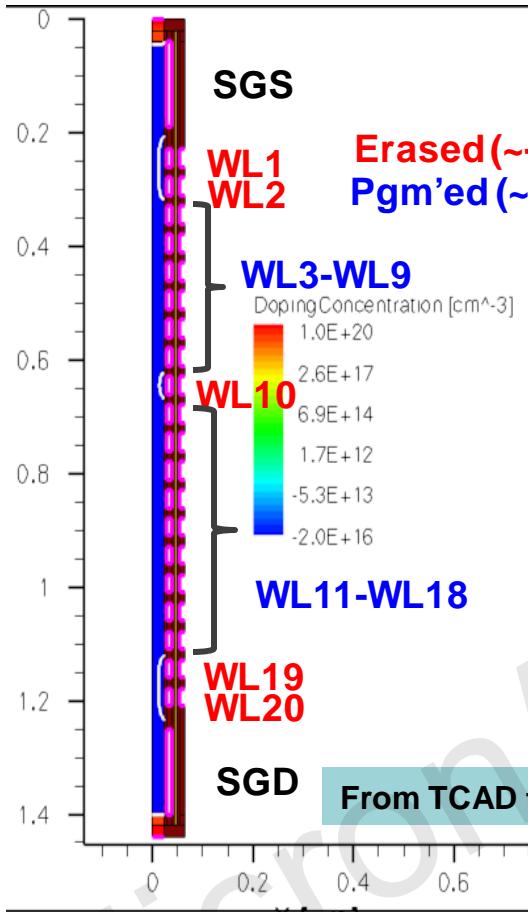
The positive potential can be sustained by the reverse biased PN junction.

- **Negative boost can't occur in 2D NAND** because the accumulation layer is connected to pwell which supplies infinite amount of holes.

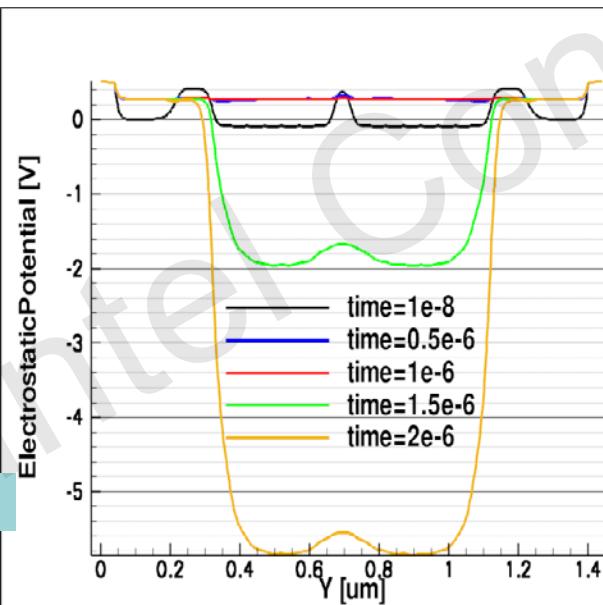
- **Negative boost can occur in 3D NAND** if the accumulation layer is sandwiched by the inversion layers. Lack of hole supply.

The negative potential can be sustained by the reverse biased virtual PN junction.

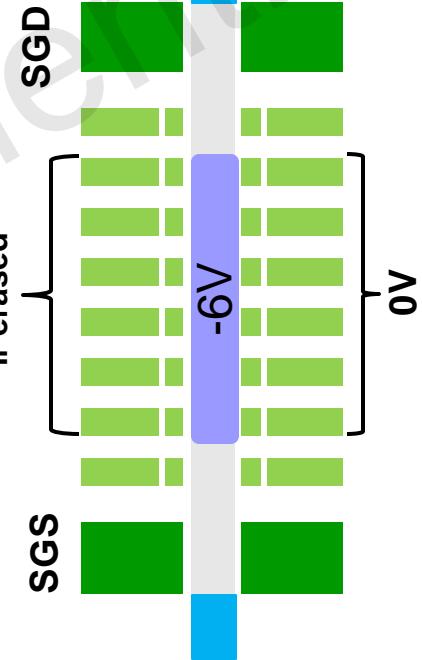
Extra read disturb stress



	Ramp up ~1us	Ramp down ~1us
WLS	0V to 8V	8V to 0V
N+ source	0V	0V
SG	3V	3V

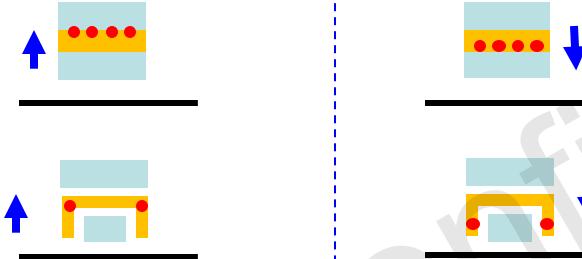


These cells sees extra read disturb stress if erased



- At the end of read after WL ramping down from 8V to 0V, the body goes to negative potential. Sustained by the virtual PN junction.
- Erased cell (WL10 in this simulation) sees excess read disturb stress due to the negative body potential.

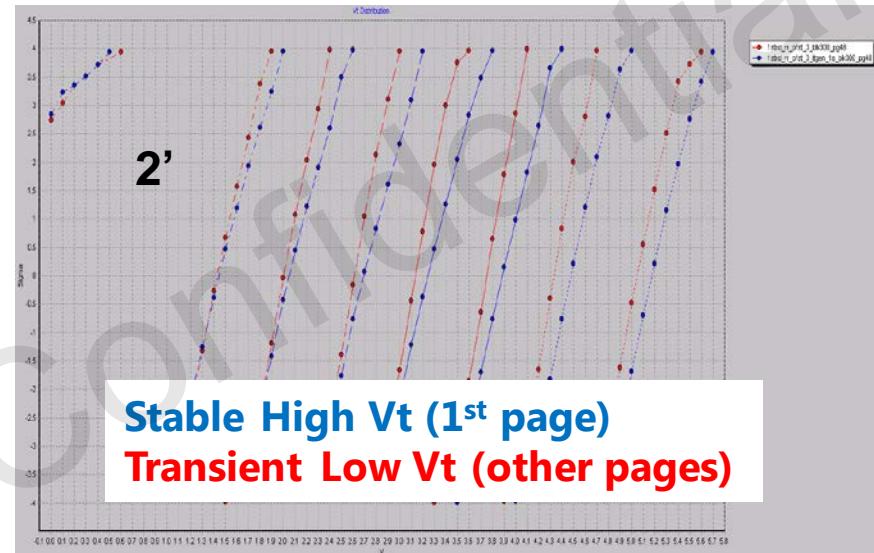
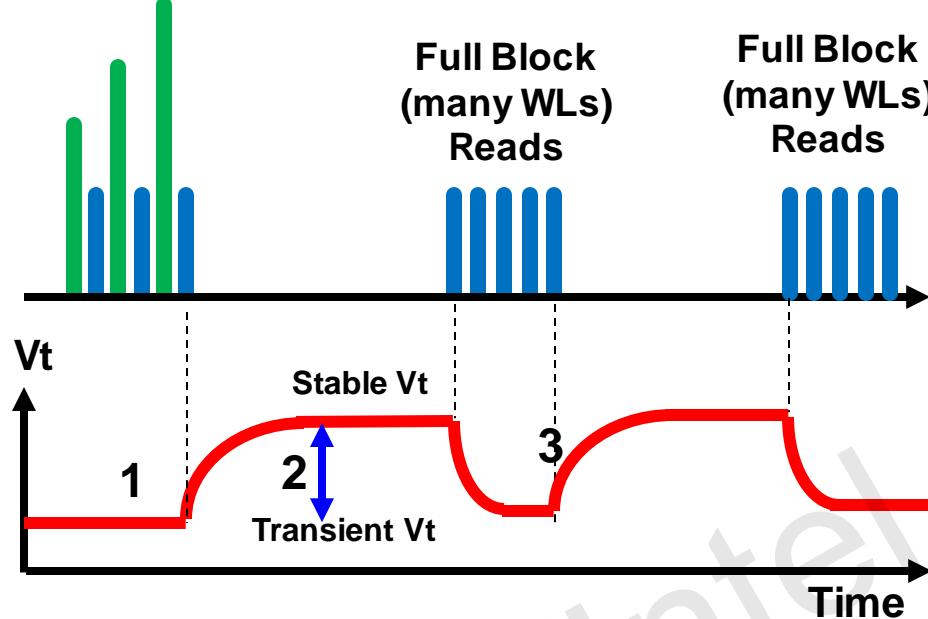
Transient Vt Mechanisms

Transient States $V_g \sim +8V$ (Read/Pgm_inh)	Stable states $V_g \sim 0V$ (idle time)	Transient Vt -> Stable Vt
Charge re-distribution or Polarization at IPD		Low VT to High VT (QCG)
Charge trapping at Ox/Si interface		High VT to Low VT (QCL)
Mobile charge diffusion		QCL when +Vg attract e. QCG when neg. body moves e away. (= cell with e- is stable)
Negative body potential	<p>$\sim +8V \rightarrow 0V$ (ramp down)</p> 	<p>Neg. body can trigger the transient.</p> <p>Or Neg. body can latch The transient states Created by +Vg.</p>

A weak gate/body voltage creates a transient Vt state without changing the # charges at the storage (FG), leading quick charge loss or gain (QCL or QCG)

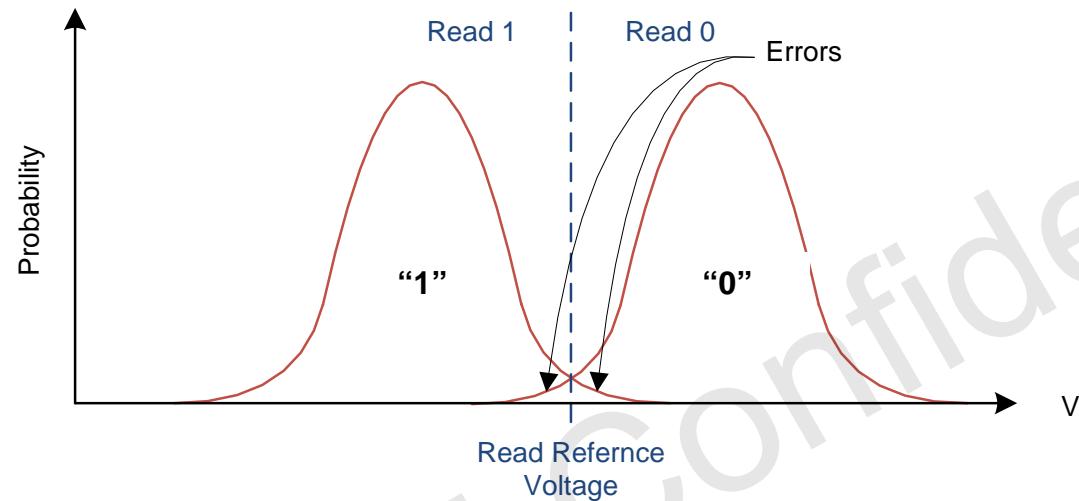
First Page Read (QCG example)

Program & Verify



1. Pgm&Verify includes gate stress, hence reading the transient V_t .
- 2 & 2'. V_t transient to the stable V_t . Delta between stable and transient V_t varies among cells, therefore, V_t distribution shifts and widens.
3. 1st page read after the idle time reads the stable V_t . As read progresses, V_t moves to the transient V_t due to V_{passR} stress.

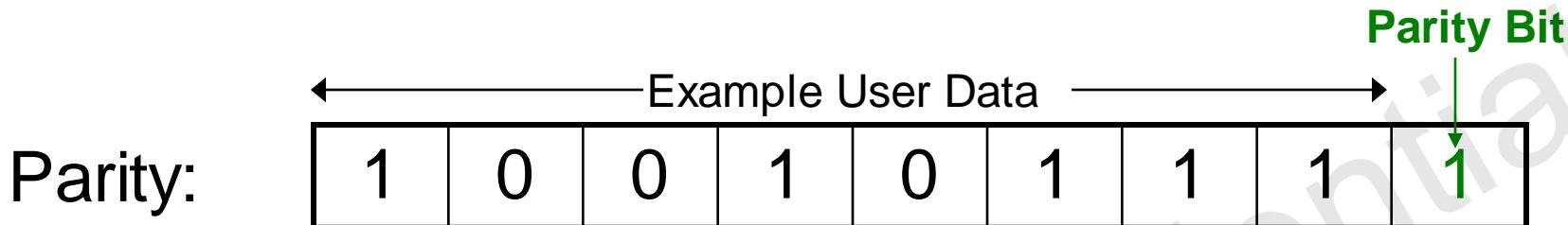
Source of Error in NAND Flash



NAND distributions overlap

- Program noise
- Erratic over programming
- Program disturb
- Read disturb
- Charge Loss (ICL, SBCL)
- X-Temp
-

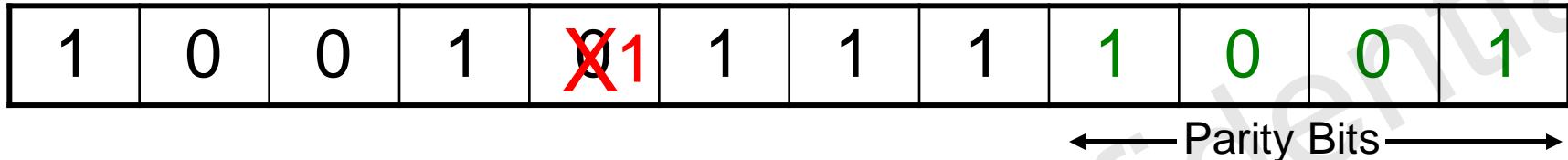
ECC: Example



- Parity is a familiar scheme for detecting single-bit errors
 - 9th bit is added to every byte and set to a value that turns the sum of all 9 bits into an even number
 - If any bit flips to the wrong state, the bits won't add up right, and the system will detect an error
 - Result: The scheme detects any single-bit error, can't correct any error, and fails if there are 2 errors
- Adding one parity bit to every byte is a scheme that localizes an error to within 9 bits

ECC: Simple Example

- By adding *four* parity bits to every byte, the error can be localized to the bad bit, and *corrected*



- Each ECC bit is a parity check on some subset of the 8 user bits:
 - 1st ECC bit = parity(bits 4,5,6,7)
 - 2nd ECC bit = parity(bits 1, 2, 3, 5, 6)
 - 3rd ECC bit = parity(bits 0, 2, 3, 5, 6)
 - 4th ECC bit = parity(bits 0, 1, 3, 4, 6)
- Example: Suppose user bit 4 fails ... how could the system figure that out?
 - Bit 4 is checked by the 1st and 4th ECC bits, which both flag the error
 - These two ECC bits share only bits 4 and 6, so the error must be one of those two
 - But bit 6 is checked by the 2nd and 3rd ECC bits, which are unaffected by the error, so bit 6 can't be it
 - Therefore, bit 4 must be the error
- The exercise will work on any of the bits (even the ECC bits themselves)
- This simple “Hamming” ECC detects and corrects a single error in a byte, but fails on >1 bit

How is Error Correction Done?

- For every k data bits, add m parity bits to generate a $n=k+m$ size codeword
- Allow errors to occur, but have parity bits which can be used to correct the errors.
- Waste of bits due to addition of m parity bits per k data bits, but huge gains in UBER (uncorrected bit error rate), for a give RBER (raw bit error rate) so worth it.
- How many parity bits do we need? Depends on the RBER of the NAND component and the level of performance (UBER) expected from NAND.
- Large code-word sizes can improve the correction efficiency
- Once we specify the RBER and the UBER, and the code-word size, parity bits required can be calculated.



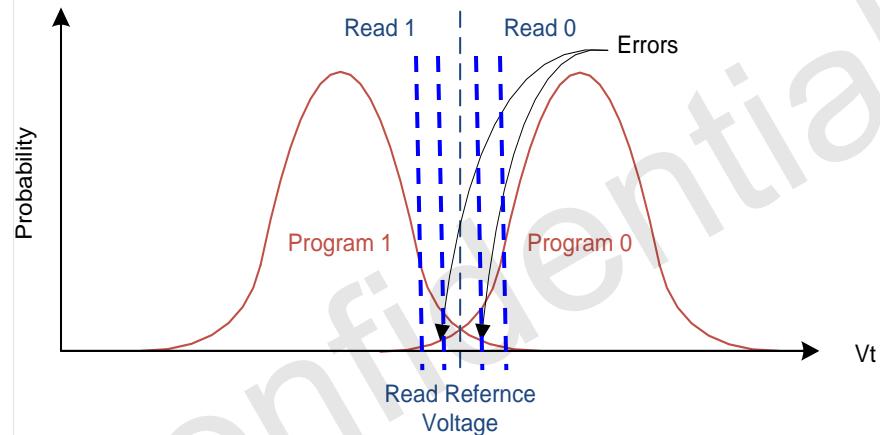
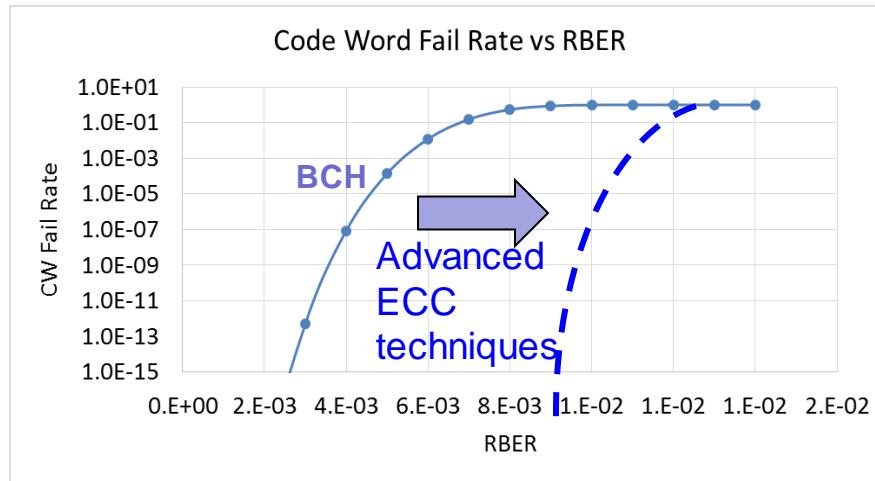
ECC Math

- Suppose one knows that a fraction P_{bit} of bits will fail in a device, called the “raw bit error rate”
- The probability that the CW will have “b” failing bits is given by the binomial probability formula:

$$P_{\text{CW}} = \frac{E!}{(E-b)!b!} \cdot P_{\text{bit}}^b \cdot (1-P_{\text{bit}})^{E-b}$$

- Here E is the # of bits in the ECC CW ($1\text{kB} + 126 \text{ bytes} = 9200 \text{ bits}$) and “!” means “factorial”
- If the ECC scheme can handle “ECC” failing bits, then the ECC fails if the # of failing bits is “ECC+1” or greater.

ECC Capability



- Based on spare bytes on our dies (16KB + 2208 ECC Bytes) we can correct up to 72 bits/1kByte correction capability.
- This doesn't mean we can tolerate 72bit/(1kB*8) of average error rate, since then every other codeword will have >72 bits failing, which will be uncorrectable.
- To guarantee good outgoing die quality the raw bit error rate (average error rate) needs to be much lower (~3.3E-3 ~27 bits/1kByte).
- Since the uncorrected error is very steep function(power of ~72), very important to make sure that worst page on the die stay below the RBER limit.
- More advanced ECC techniques that can use analog Vt information are used to improve correction capability and allow for even higher RBER (up to ~1E-2).

~Appendix~

Impact of Trapped Charge on Cell Program/Erase

Charge Trapped in Tunox - UV_Vt, VgVt, Vtp_sat

- Consider Charge Trapped at a distance “d” from the Tunox-Si interface. (Assume $d > 25\text{A}$ to ignore barrier lowering and impact to tunneling)
- EOT is $(\text{Tox} + \text{IPD}_{\text{eot}}/\text{Wrap}_{\text{Ratio}})$ or Tox/GCR .

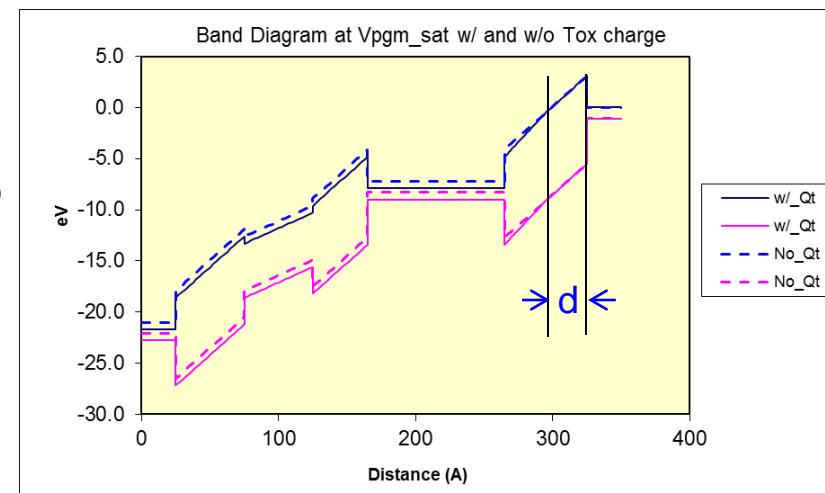
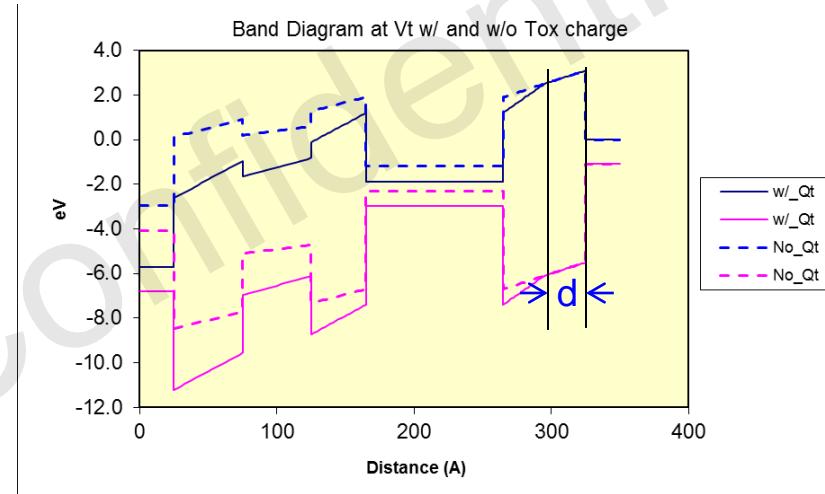
At $V_g=V_t$ or V_g corresponding to onset of programming ($E_{\text{ox}}=E_{\text{pgm}} (\sim 12\text{MV/cm})$), the E_{field} in the tunnel-ox at Ox/Si interface is unchanged by trapped charge. But E_{field} above the trapped charge is higher by (Q/ϵ) .

- $\Delta_{\text{UV}_Vt} = (-Q/\epsilon) * (\text{EOT}-d)$
- $\Delta_{\text{Vpgm}} = (-Q/\epsilon) * (\text{EOT}-d)$
- $\Delta_{\text{VgVt}} = 0$

(If trapped charge is electrons, then $-Q$ becomes a +ve quantity and hence UV_Vt shifts in the +ve direction)

Again at Vpgm_{sat} , the E_{field} in IPD is same w/ or w/o Tox trapped charge ($\sim 12\text{MV/cm}$). Only region that has higher E_{field} is region between “d” and “FG”

- $\Delta_{\text{Vpgm}_{\text{sat}}} = (-Q/\epsilon) * (\text{Tox}-d)$
- $\Delta_{\text{Vtp}_{\text{sat}}} = (-Q/\epsilon) * (\text{Tox}-d)$



Charge Trapped in Tunox – VwVt, Vte_sat

- Consider Charge Trapped at a distance “d” from the Tunox-Si interface. (Assume $d > 25\text{A}$ to ignore barrier lowering and impact to tunneling)

For Erase, we need to maintain the same E-field (Erase $\sim 12\text{MV/cm}$) at the FG/Tox interface. Which means same E-field in the IPD. The only region with higher E-field is the region between “d” and Si/Ox interface.

➤ $\Delta_{\text{Vera}} \text{ (Well voltage)} = (-Q/\epsilon) * (d)$

(**Vera here refers to the +ve Well voltage**)

And we had:

➤ $\Delta_{\text{UV}} \text{ Vt} = (-Q/\epsilon) * (\text{EOT} - d)$

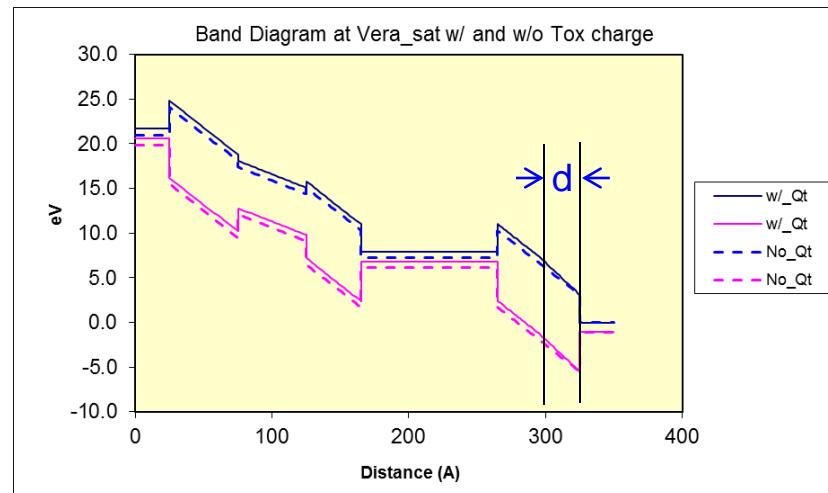
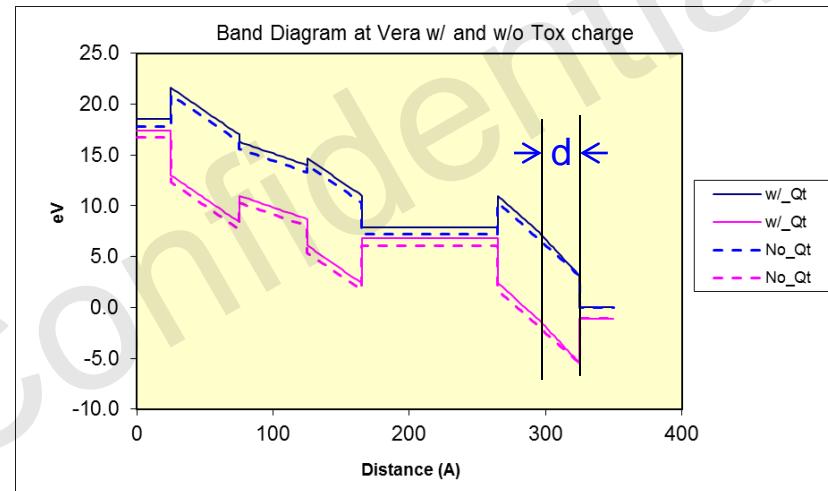
➤ $\Delta_{\text{Vw}} \text{ Vt} = (-Q/\epsilon) * (\text{EOT})$

- At Vera_sat, the E-field in IPD is same w/ or w/o Tox trapped charge ($\sim 12\text{MV/cm}$). Only region that has higher E-field is region between “d” and Si/Ox interface.

➤ $\Delta_{\text{Vera}} \text{ sat} = (-Q/\epsilon) * (d)$

➤ $\Delta_{\text{Vte}} \text{ sat} = (-Q/\epsilon) * (\text{EOT} - d)$

(**Vte_sat here is the actual Vt (NOT magnitude) = VwVt – Vera_sat**)



Charge Trapped in IPD - UV_Vt, VgVt, Vtp_sat

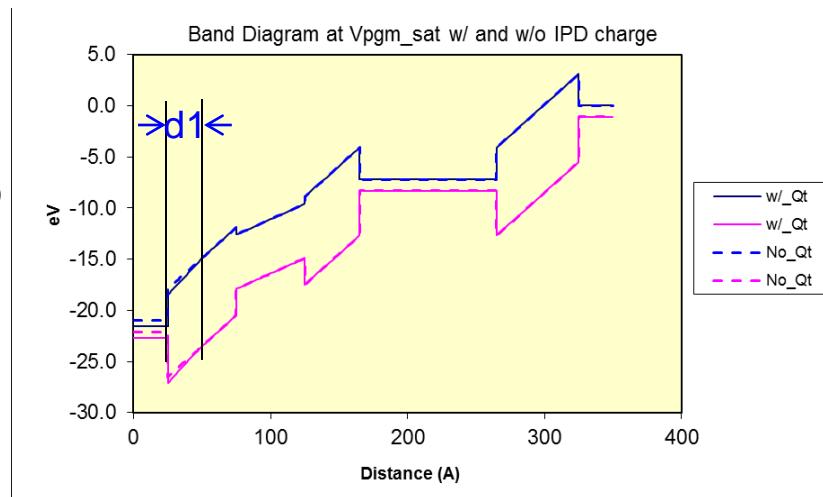
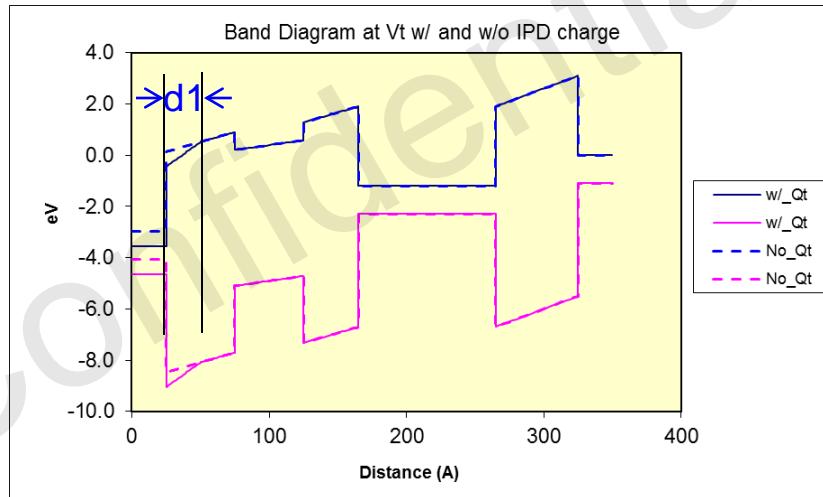
- Consider Charge Trapped at a distance “d1” from the CG-IPD interface. (Assume $d > 25\text{A}$ to ignore barrier lowering and impact to tunneling).
- Even though schematically the charge is shown trapped in the top IPD, the formula is applicable charge trapped in Nitride or bottom IPD as well.

At $V_g = V_t$ or V_g corresponding to onset of programming ($E_{ox} = E_{pgm}$ ($\sim 12\text{MV/cm}$)), the E_{field} in the tunnel-ox at Ox/Si interface is unchanged by trapped charge. But E_{field} above the trapped charge is higher by (Q/ϵ) .

- $\Delta_{UV_Vt} = (-Q/\epsilon) * (d_1)$
- $\Delta_{V_{pgm}} = (-Q/\epsilon) * (d_1)$
- $\Delta_{VgVt} = 0$

Again at V_{pgm_sat} , the E_{field} in IPD is same w/ or w/o trapped charge ($\sim 12\text{MV/cm}$). Only region that has higher E_{field} is region between “ d_1 ” and “CG”

- $\Delta_{V_{pgm_sat}} = (-Q/\epsilon) * (d_1)$
- $\Delta_{V_{tp_sat}} = (-Q/\epsilon) * (d_1)$



Charge Trapped in IPD – VwVt, Vte_sat

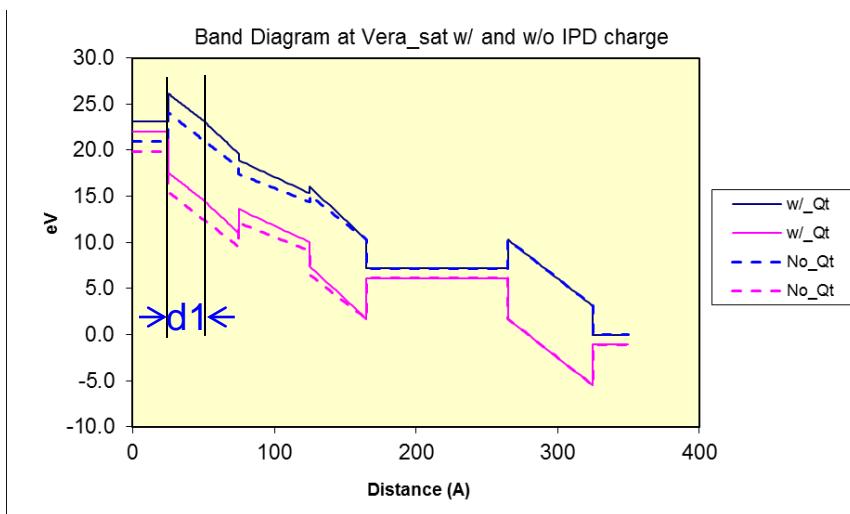
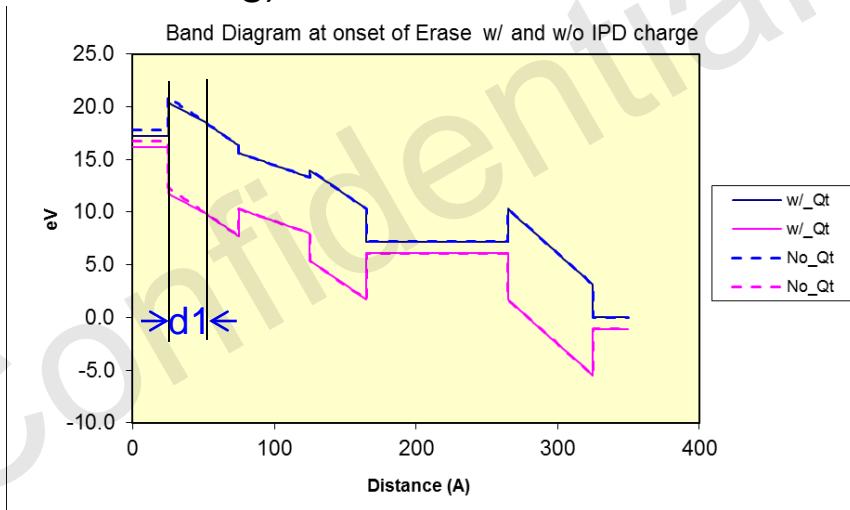
- Consider Charge Trapped at a distance “d1” from the CG-IPD interface. (Assume d >25A to ignore barrier lowering and impact to tunneling)

For Erase, we need to maintain the same E-field (Erase ~12MV/cm) at the FG/Tox interface. Which means same E-field in the IPD. The only region with lower E-field is the region between “d1” and CG/IPD interface.

➤ $\Delta_{Vera} = (Q/\epsilon) * (d1)$
(Vera here refers to the +ve Well voltage)

And we had:

- $\Delta_{UV_Vt} = (-Q/\epsilon) * (d1)$
- $\Delta_{Vw_Vt} = 0$
- At Vera_sat, the E-field in the top part of IPD is same w/ or w/o IPD trapped charge (~12MV/cm). But the region below d1 through FG will have higher E-field.
- $\Delta_{Vera_sat} = (-Q/\epsilon) * (IPD_EOT - d1)$
- $\Delta_{Vte_sat} = (Q/\epsilon) * (IPD_EOT - d1)$
(Vte_sat here is the actual Vt (NOT magnitude) = VwVt – Vera_sat)



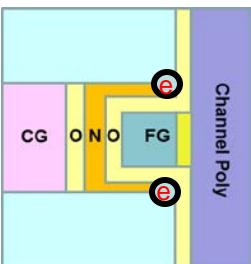
Calculator

		Tox	IPD3	IPD1	Comment
-Qt/ε	V/cm	2.32E+06	2.32E+06	2.32E+06	Charge density/cm^2
Tox	A	60	60	60	Tox thickness
IPD	A	115	115	115	IPD EOT
Wrap Ratio	#	1.3	1.3	1.3	Plecholder wrap ratio
EOT	A	148.5	148.5	148.5	Tox + (IPD EOT / Wrap Ratio)
d	A	30			Tox Trap location from Channe linterface
d1	A		95.0	25.0	IPD Trap location from CG interface
					Tox Eq
$\Delta_{UV_Vt} =$	V	2.74	2.20	0.58	$\Delta_{UV_Vt} = (-Q/\varepsilon)^*(EOT-d)$
$\Delta_{VgVt} =$	V	0.00	0.00	0.00	$\Delta_{VgVt} = 0$
$\Delta_{VwVt} =$	V	3.44	0.00	0.00	$\Delta_{VwVt} = (-Q/\varepsilon)^*(EOT)$
$\Delta_{Vpgm_sat} =$	V	0.70	2.20	0.58	$\Delta_{Vpgm_sat} = (-Q/\varepsilon)^*(Tox-d)$
$\Delta_{Vera_sat} =$	V	0.70	0.46	2.09	$\Delta_{Vera_sat} = (-Q/\varepsilon)^*(d)$
$\Delta_{Vtp_sat} =$	V	0.70	2.20	0.58	$\Delta_{Vtp_sat} = (-Q/\varepsilon)^*(Tox-d)$
$\Delta_{Vte_sat} =$	V	2.74	-0.46	-2.09	$\Delta_{Vte_sat} = (-Q/\varepsilon)^*(EOT-d)$
$\Delta_{P/E\ Window}$	V	-2.05	2.66	2.66	
					IPD3 / IPD1 Eq

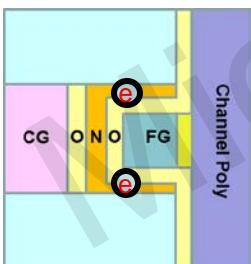
- Charge trapping in Tunox reduces P/E window.
- Charge trapping in IPD can significantly improve the P/E_sat window.

Other Charge Trapping Locations

- In addition to the Tox and IPD (directly in the path of channel V_t and/or Tunneling), we could have charge trapped around the cell (flank nitride, over LDD (beak region)).
- These charge can couple to the channel and/or FG. (But they themselves do not directly alter the Tunox or IPD E-field):
 - The coupling to channel (LDD, beak, etc.) can be thought of as UV_V_t change by channel doping change. It will impact the P/E V_{ts} the way channel doping would have impacted P/E v_{ts}. So, V_{gVt}, V_{wVt}, V_{tp_sat}, V_{te_sat}, etc. will change by the same amount as Δ_{UV_Vt} .
 - $\Delta V_t > 0$ during the read operation;
 - $\Delta V_{pgm} = 0$; $\Delta V_g V_t = \Delta V_{pgm} - \Delta V_t = -\Delta Vt < 0$; $\Delta V_{tp_sat} = \Delta Vt > 0$;
 - $\Delta V_{era} = 0$; $\Delta V_w V_t = \Delta V_{era} + \Delta V_t = \Delta Vt > 0$; $\Delta V_{te_sat} = \Delta Vt > 0$;
 - $\Delta P/E_{window} = 0$
 - The coupling to FG is really going to be indistinguishable from those charges being on the FG itself



- $\Delta V_t > 0$ during the read operation;
- $\Delta V_{pgm} = \Delta V_t$; $\Delta V_g V_t = 0$; $\Delta V_{tp_sat} = 0$;
- $\Delta V_{era} = -\Delta V_t$; $\Delta V_w V_t = 0$; $\Delta V_{te_sat} = 0$;



Charge Trapping Summary

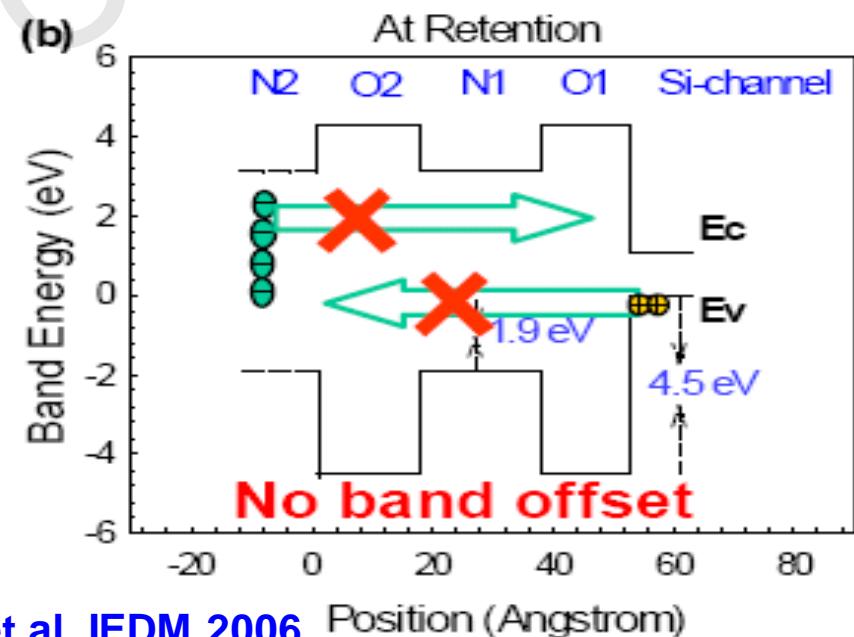
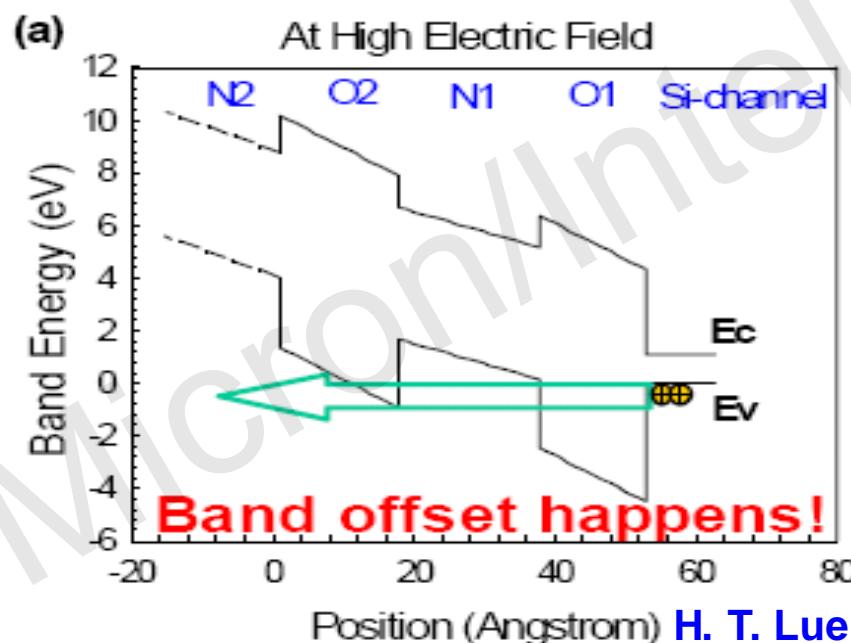
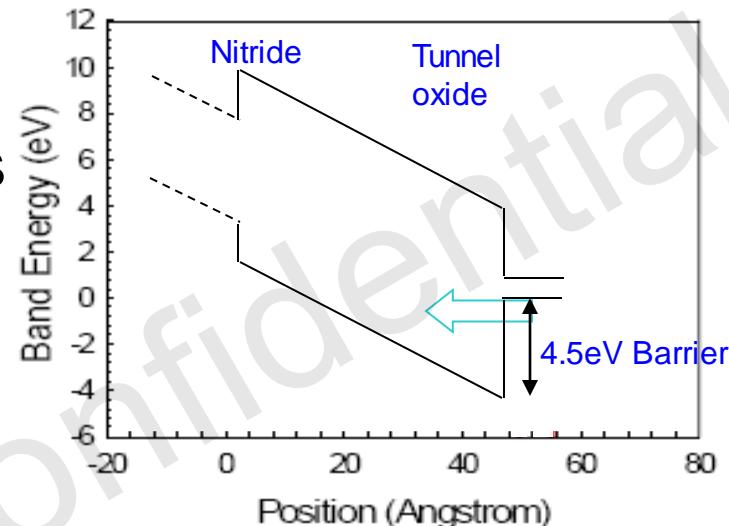
- Even though we looked at the impact of trapped charge in each of the different locations in isolation...
- ...In practice, we could have charge trapped in multiple locations simultaneously
- And also based on the location of the charge, it will couple to both the channel as well as the FG and may partially be in the path of the tunneling (Tunox or IPD) as well.
- So, the real impact to the cell P/E characteristic will be some superimposition of the various impacts described.

Back-Up

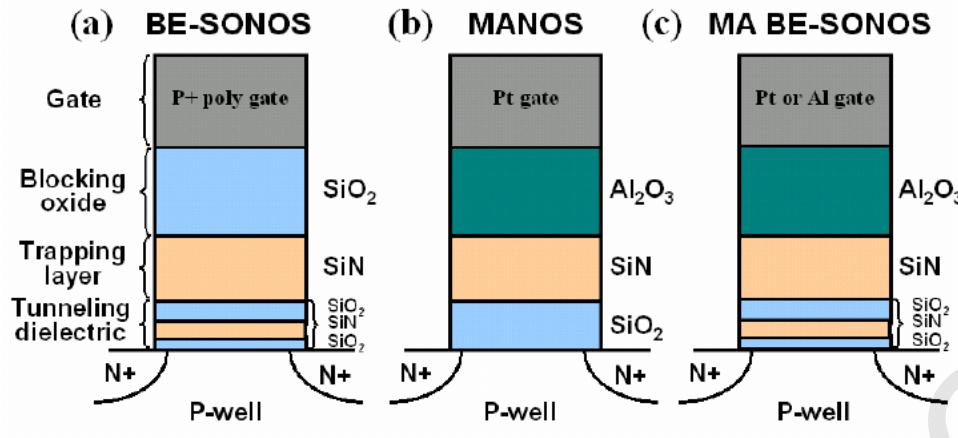
Micron/Intel Confidential

Nitride Storage - Erase

- Nitride Storage memories erase through hole injection which requires very high field.
- Engineered Tunnel barriers addresses this problem



Nitride Storage - Erase



S-C. Lai, et al. NVSMW 2007

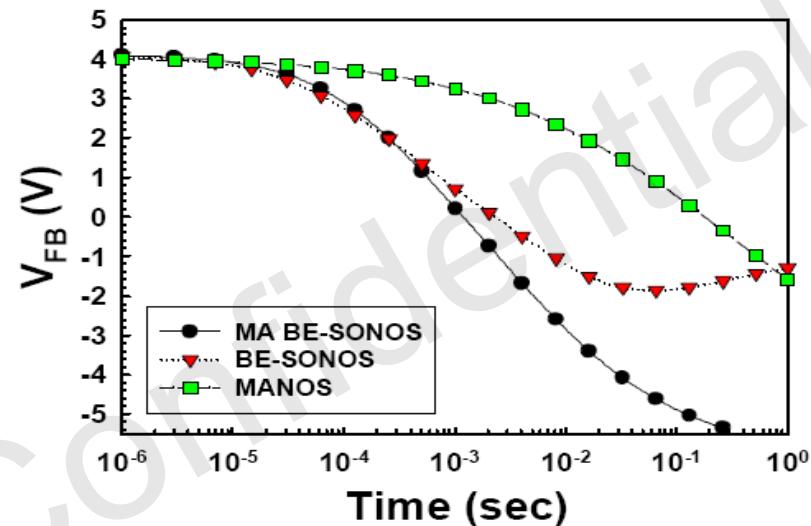
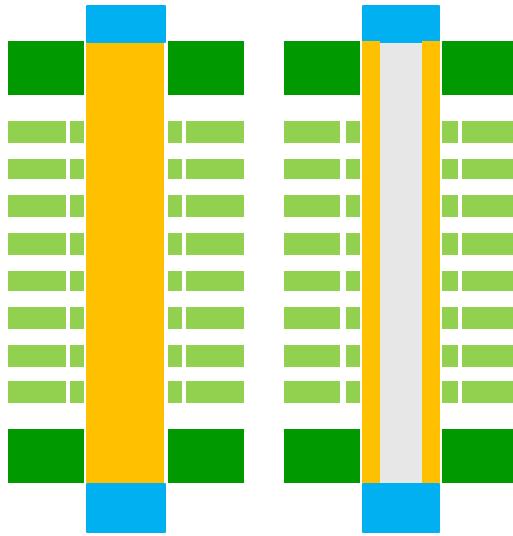


Fig. 2 Erase curves (V_{FB} -time) for the 3 devices at $V_G = -18V$, where the gate material for MA BE-SONOS and MANOS is platinum.

- Key requirements for good erase characteristics:
 - Engineered tunnel barrier
 - High work function Control Gate
 - High-K blocking dielectric

Hollow Channel

Solid channel



Hollow channel

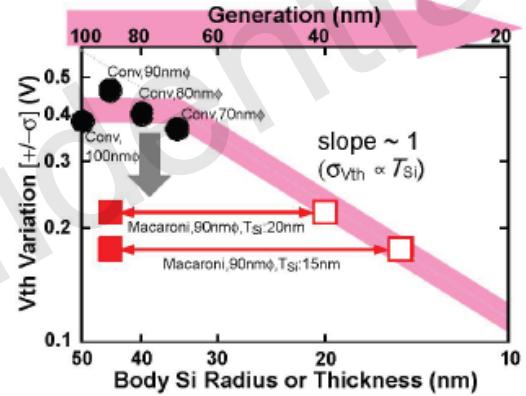
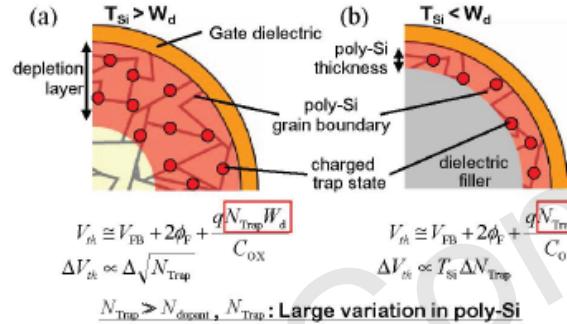


Figure 15: Extendability of ‘Macaroni’ body vertical FET. V_{th} variations in 300mm wafer are plotted as a function of body Si radius or thickness.

Fukuzumi et al., IEDM2007

- Hollow channel has three advantages due to the thin channel thickness

{ ~30nm (solid) to ~10nm (hollow) }

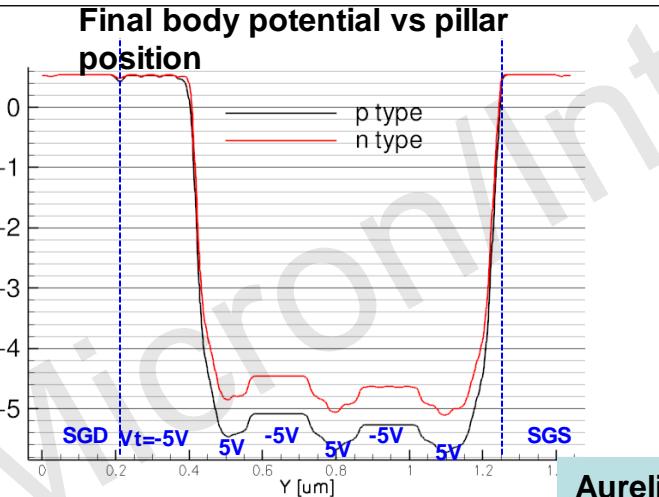
- V_t sigma (PVS) improvement due to fewer # of traps involved in the depletion layer.
- Boosting leakage reduction due to fewer G/R centers.
- Sub-threshold slope and cut-off improvement due to better channel controllability, allowing n-type channel.

Floating Body Effect - Read

- At the end of the Read waveform, the pillar body can go negative (WL ramp down)
- The resulting negative voltage can linger for some time resulting in read disturb.
- Can happen in other waveforms and contexts as well.

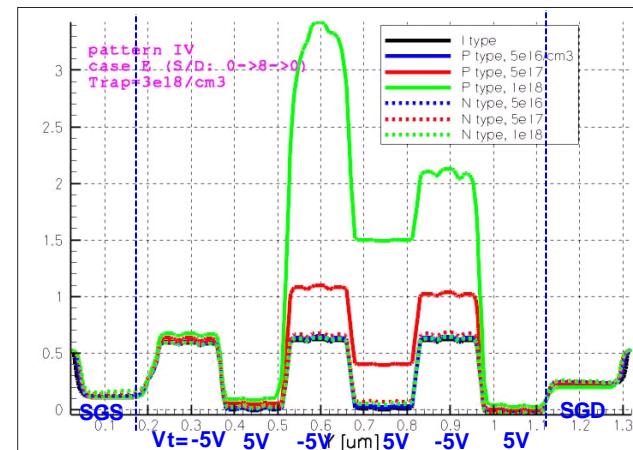
Ramp up and down is 1us

SG	WLs	N+
3V	0V to 8V to 0V	0V



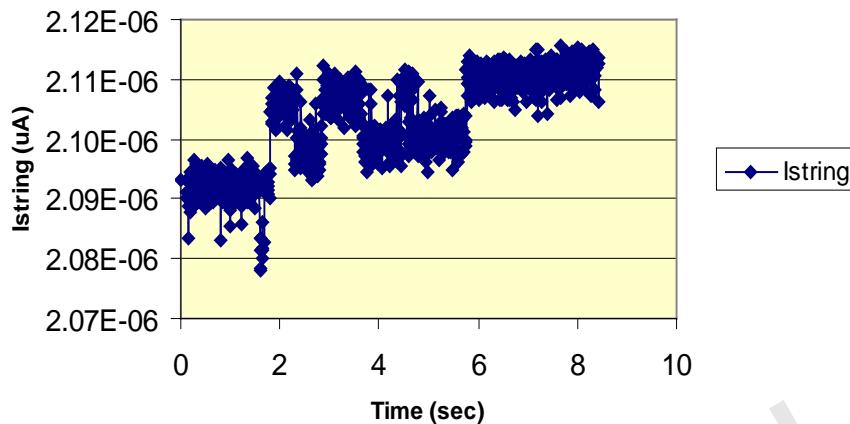
Ramp up and down is 1us

SG	N+ S/D	WL
0V	0V to 8V to 0V	0V

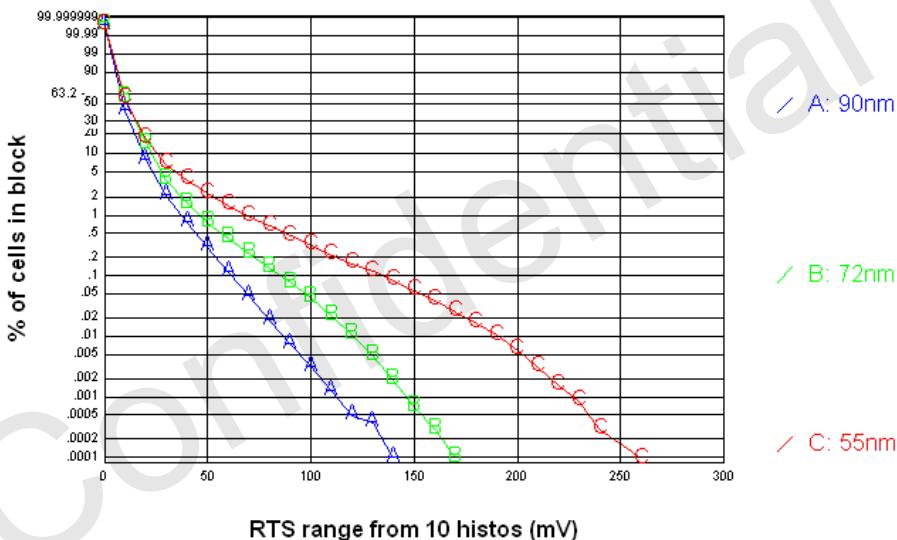


Random Telegraph Signal (RTS)

Cell Current Fluctuations in Time (Random Telegraph Signal) 55nm Node



RTS range on 90/72/55nm



- Channel electron trapping and de-trapping during read causes current to fluctuate in discrete steps, like a random telegraph key being hit -> RTS noise
- The longer you look (many reads), the bigger the events you'll see...
- The more number of cells you look, the bigger the events you'll see