



LUND
UNIVERSITY

EITP25 2020

Lecture 10 – Magnetic Memories

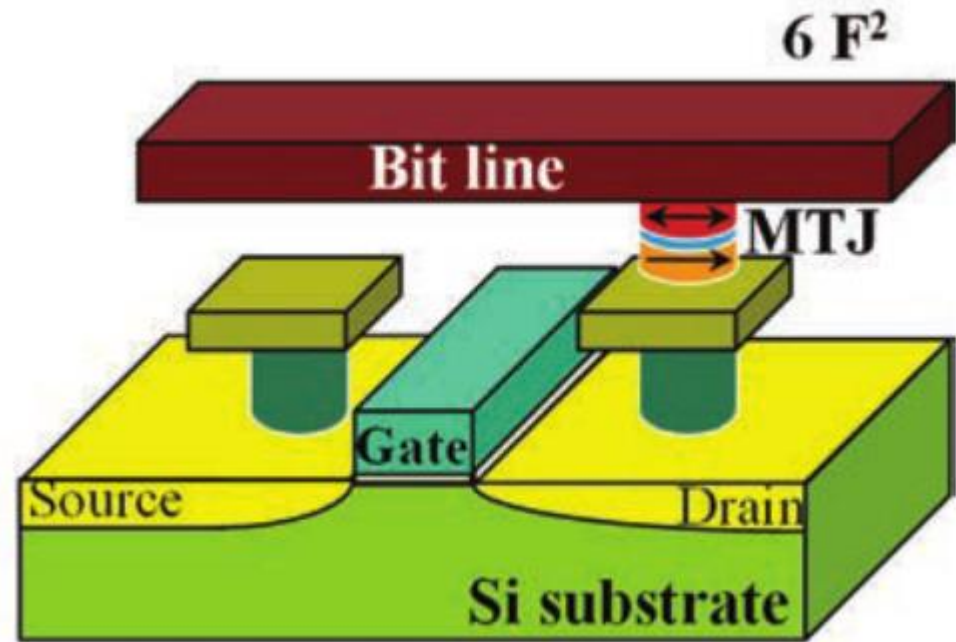


Outline – Lecture 10

- Magnetic memory principle
- The magnetic tunnel junction
- Spin-torque transfer
- Scaling STT-MRAM
- Stoichastic synaptic devices

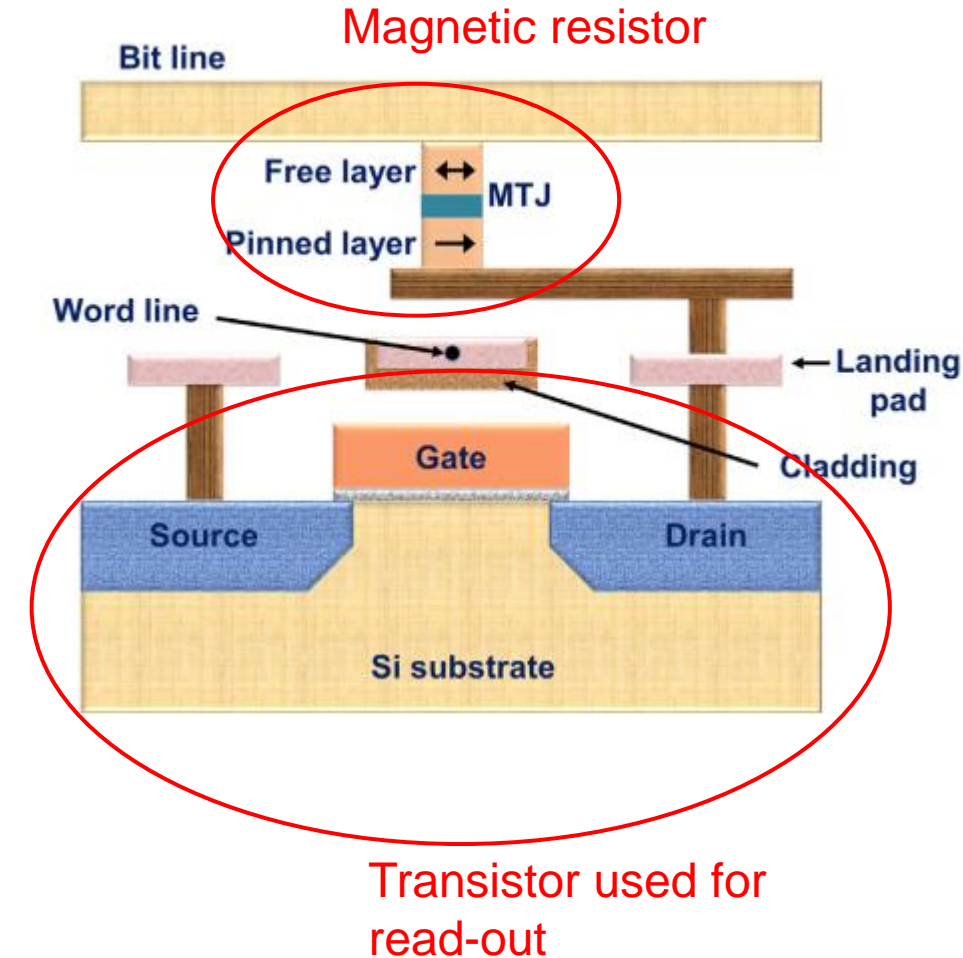
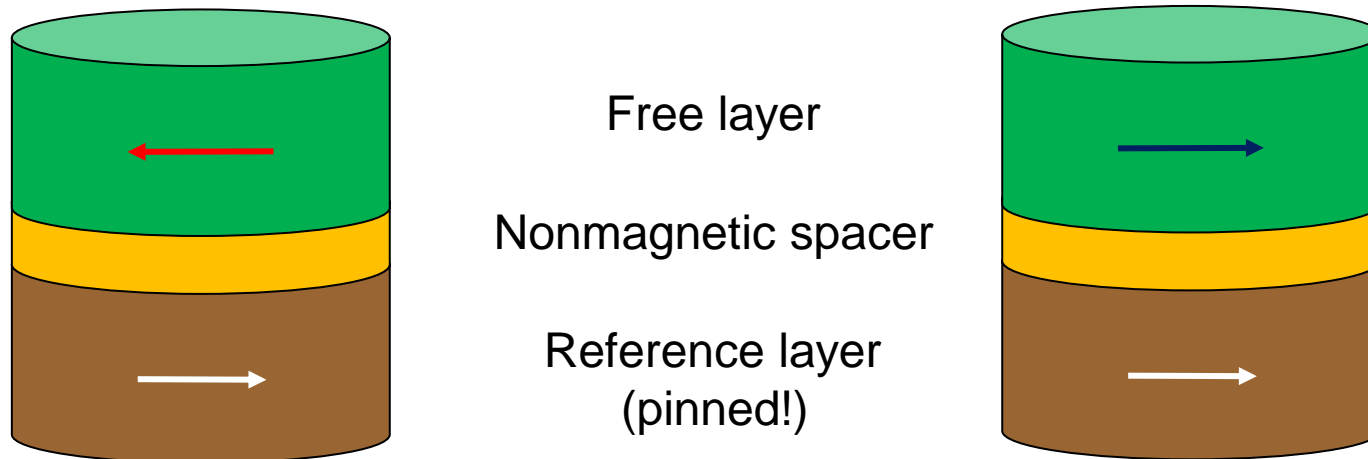
Promise of MRAM

- Energy
 - SRAM: 10's fJ/bit
 - DRAM: 1-10 pJ/bit
 - **STT-MRAM: 10-100 fJ/bit**
- Speed/delay
 - SRAM: 1 ns
 - DRAM: ~10 ns
 - **STT-MRAM: 1-10 ns**
- Cost/Density
 - SRAM: $120F^2$
 - DRAM: $6F^2$
 - **STT-MRAM: $4-6F^2$**
- Volatility
 - SRAM and DRAM: Volatile
 - **STT-MRAM: Non-volatile**



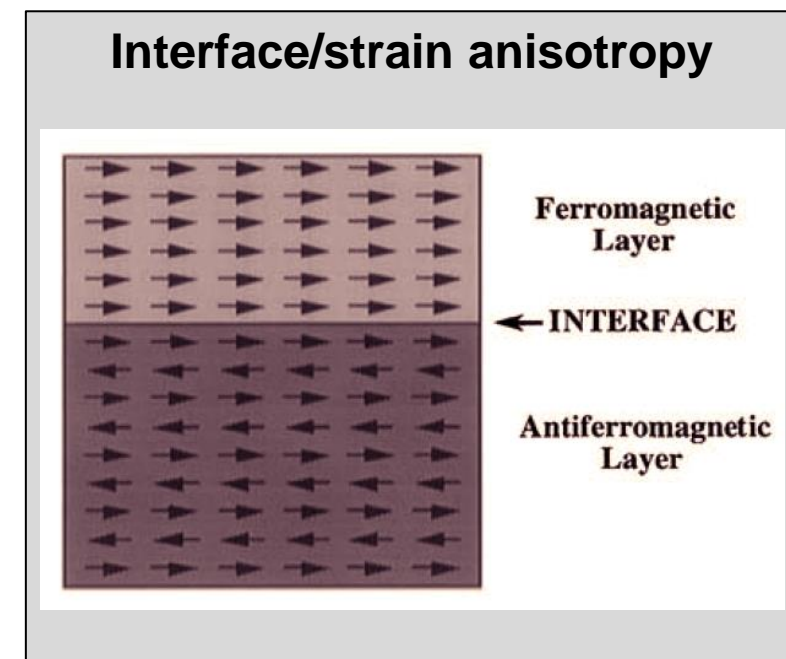
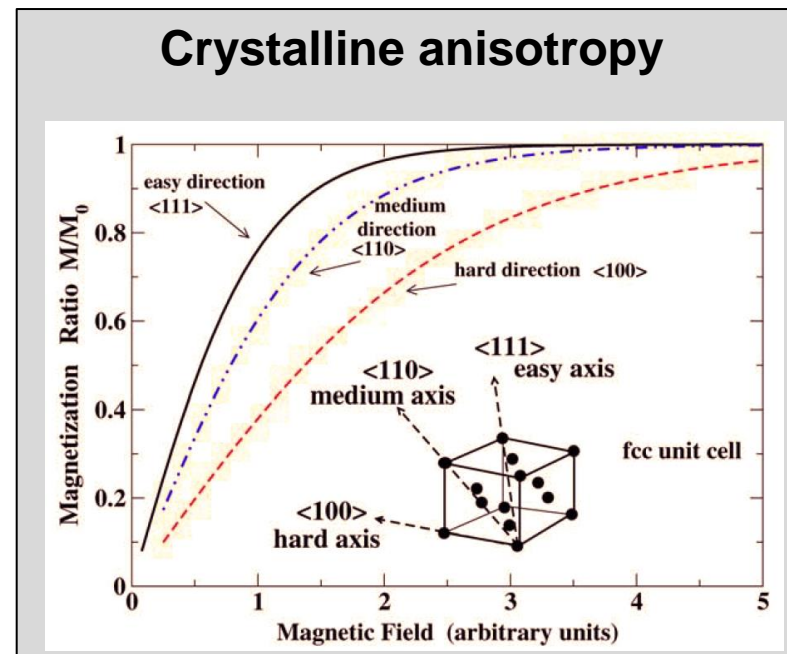
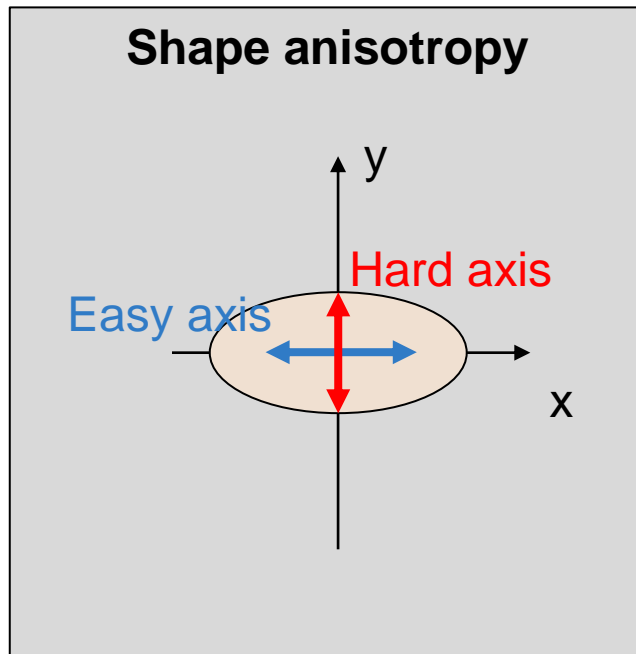
Principle of Magnetic RAM

- 1T1R device → Minimum size $6F^2$
- Memory state is stored in magnetization of a thin film
- Magnetic element connected in series with transistor
- Read-out by difference in resistance between the two magnetization states of the free layer



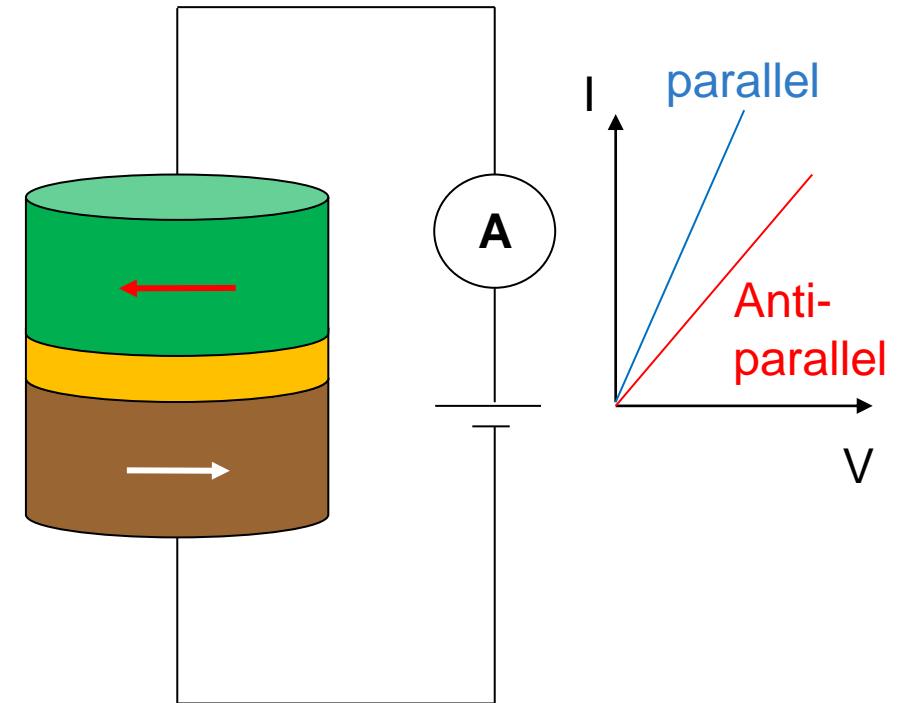
Magnetic Anisotropy

- Contrary to *magnetic isotropy*, magnetization is preferable in certain directions.
- Easy axis: An axis in which magnetization is preferable



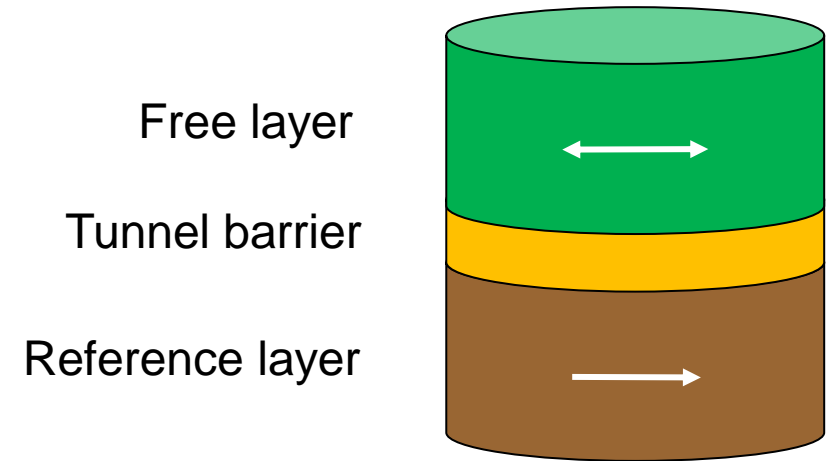
Reading out state - Magnetoresistance

- First observed in 1856, in Fe and Ni
- Due to spin-orbit coupling: Magnetization perturbs the electron cloud \rightarrow changes amount of scattering
- Magnitude $\sim 2\%$
- Giant Magnetoresistance (discovered in 80's) in Cr/Fe stacks
 - Spin-dependent scattering
 - \rightarrow Harddrive revolution
 - Albert Fert and Peter Grünberg Nobel Prize in Physics 2007.
 - Magnitude $\sim 10\text{-}20\%$



Magnetic Tunnel Junction (MTJ)

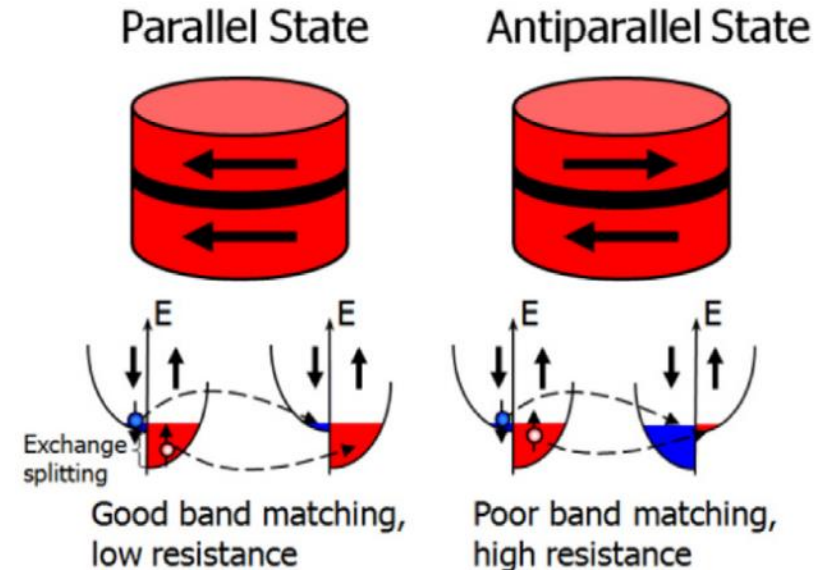
- Still Free layer and Reference layer
- Spacer → NOT conductive
- Transport through quantum tunneling
 - Tunnel barrier typically Aluminum oxide



→ Tunneling Magnetoresistance (TMR)

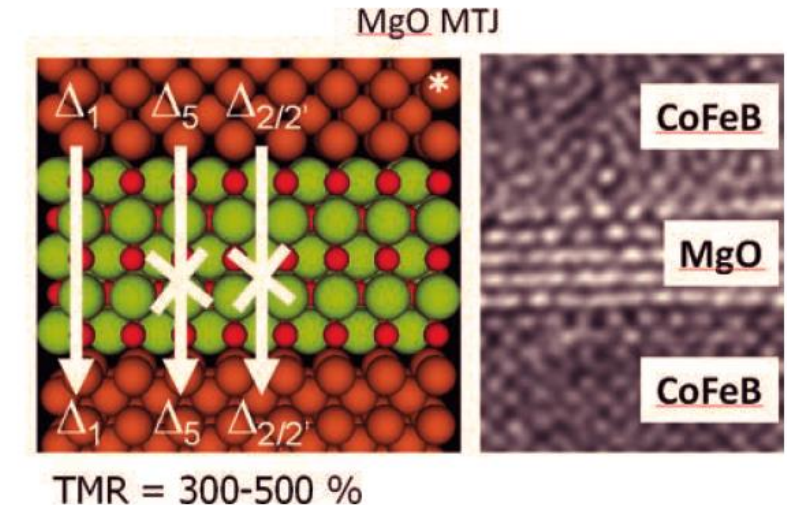
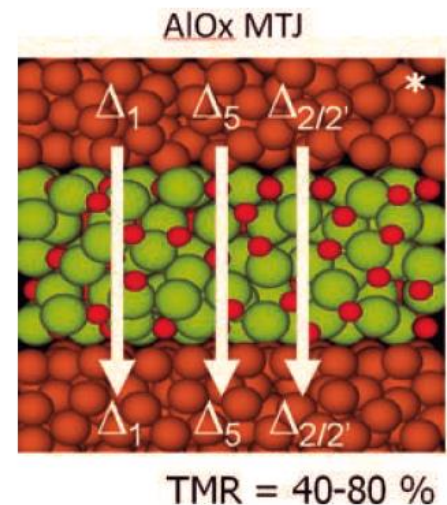
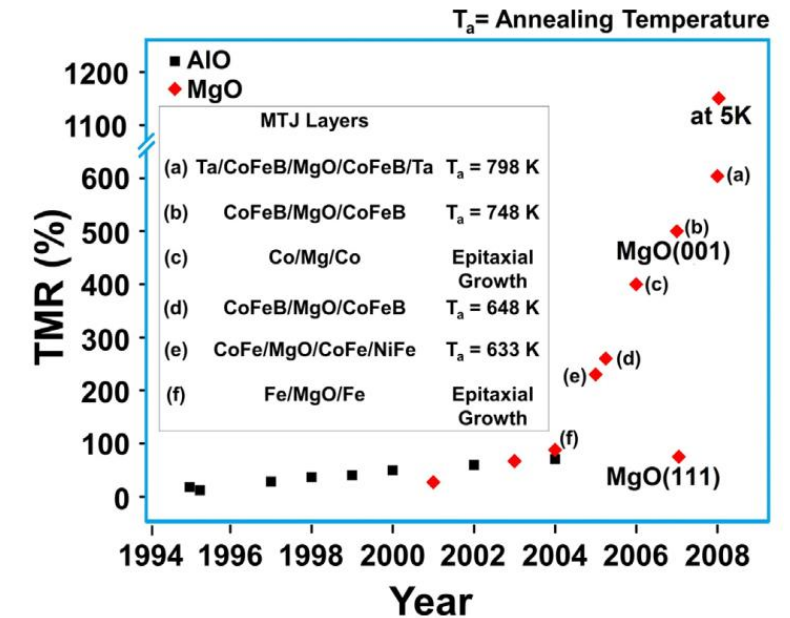
$$TMR = \frac{R_{ap} - R_p}{R_p} \sim 70 \%$$

Alignment of spin gives higher tunneling probability



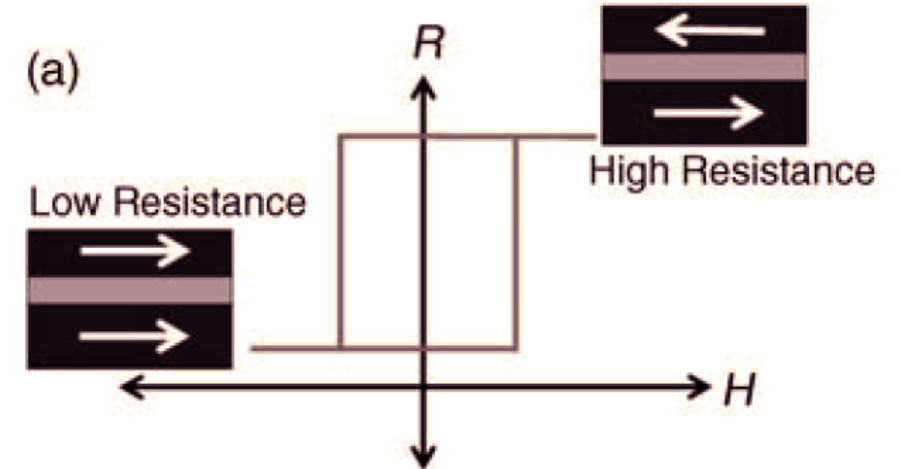
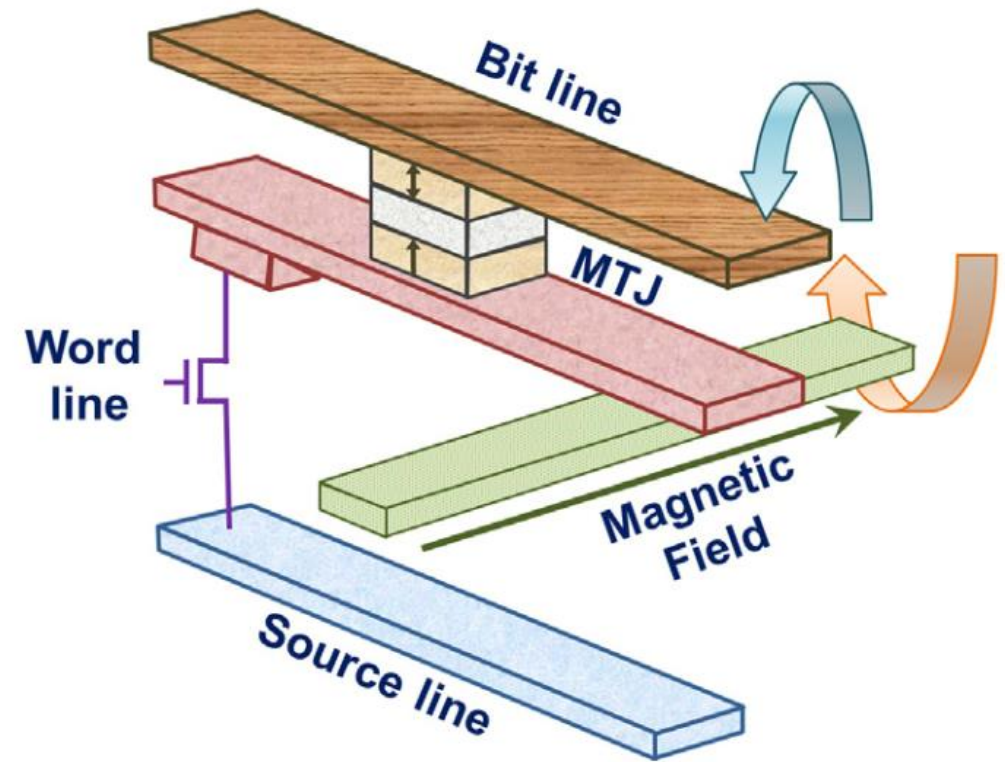
"Giant" Tunneling magnetoresistance

- **Giant TMR through crystalline MgO barriers**
 - spd hybridized states (Δ_1) are fully polarized
 - Couple to evanescent Δ_1 states in MgO
 - MgO(001) evanescent Δ_1 decay slowest
 \rightarrow filters out states with low spin polarization
 \rightarrow Makes TMR effect very strong!
 - TMR ~ 600%



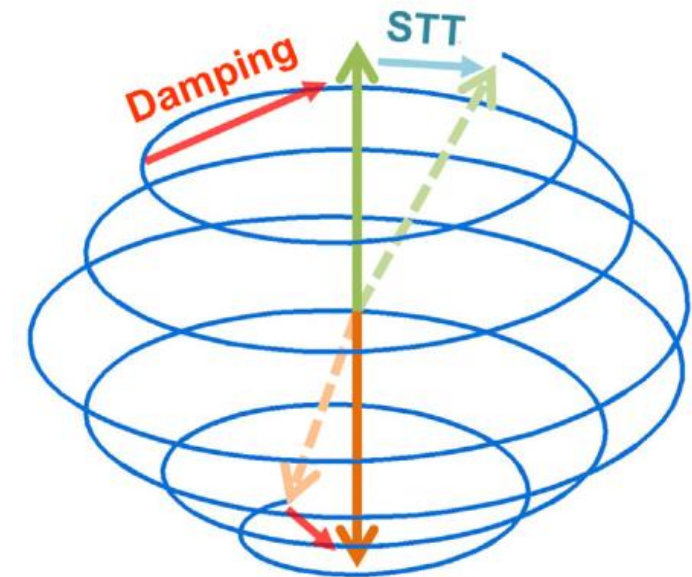
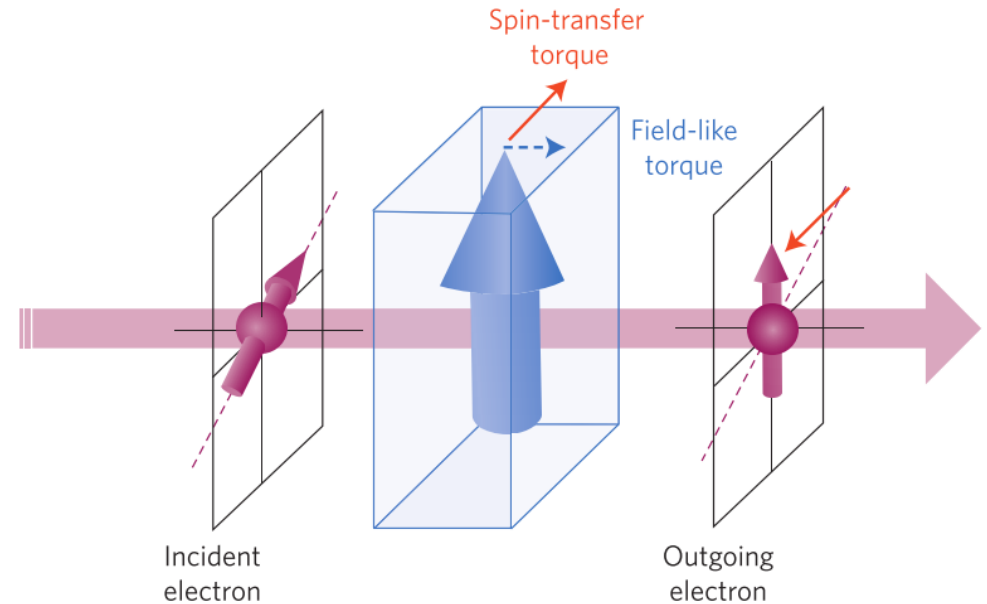
Writing in MRAM

- Current pulses on BIT and WORD lines create orthogonal B-fields
- Only sum of these magnetic field strengths should be enough to switch magnetization in Free Layer.
→ Selection of device possible
- Current needed for sufficient field strengths $> H_K^{2/3}$ prevents scaling as $H_K \propto \frac{1}{L}$, where L is device size
- Needed other way to operate MRAM...



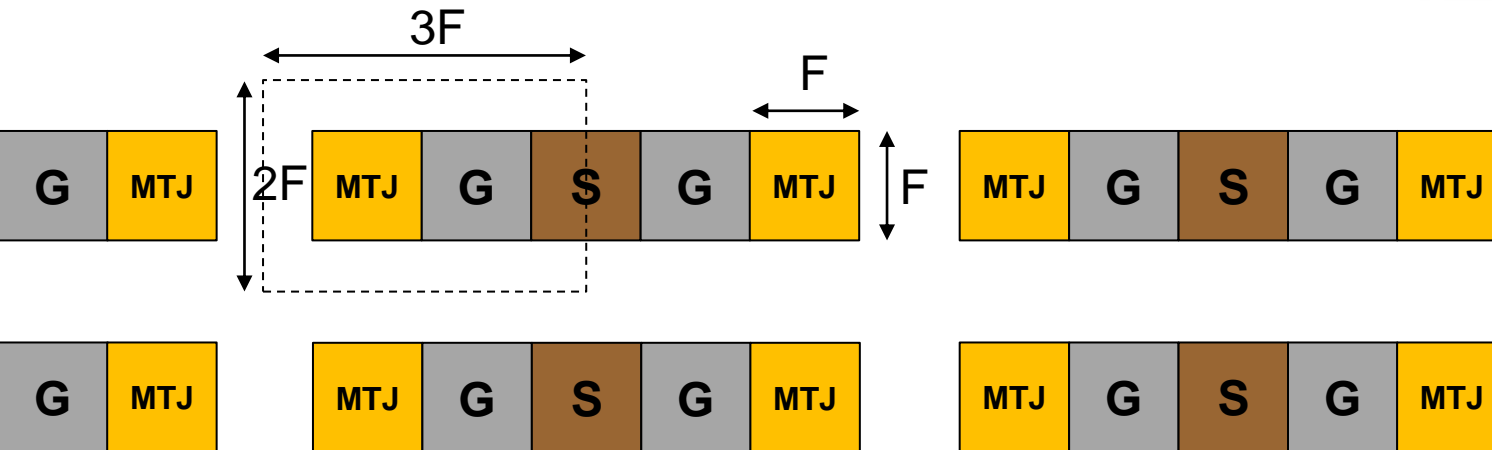
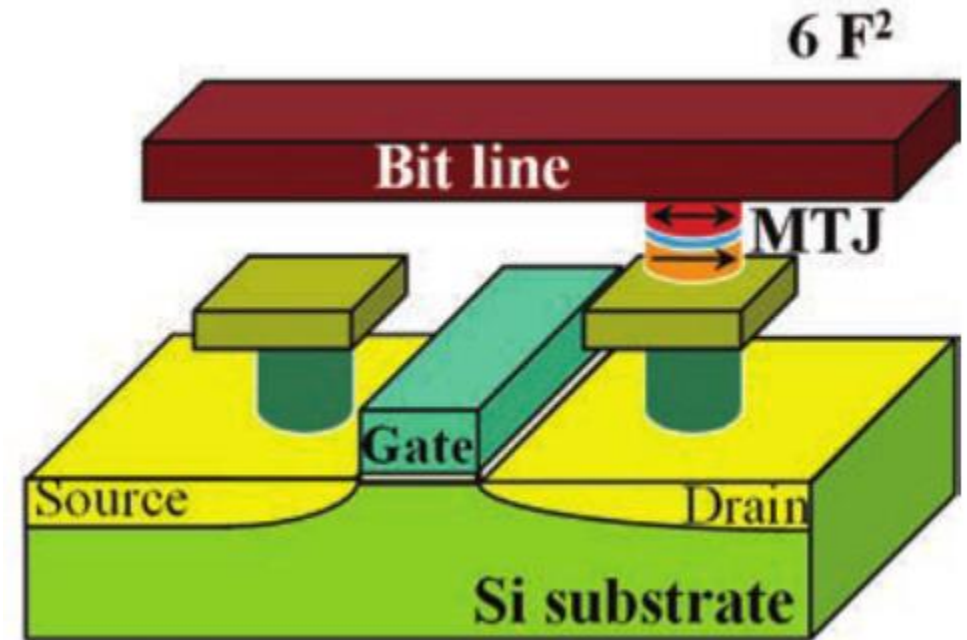
Spin-torque transfer

- Electrons flowing through the tunnel junction can transfer spin angular momentum
 - Resulting "torque" on magnetization
 - **WRITING** of data possible
-
- An electron spin passing through a magnetized layer has its spin direction altered
 - But also the magnetization is affected (torque)
 - Magnetization goes through a precession path with *damping*



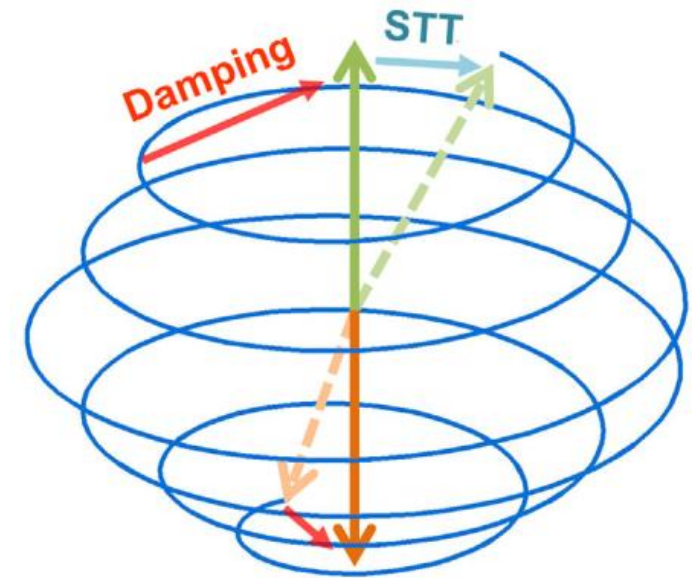
STT-MRAM layout

- Transistor Gate \rightarrow Word line
- MTJ on Drain \rightarrow Bitline
- No external magnetic field needed!
- Switching magnetization by current (STT)
- Minimum memory cell area: $6F^2$



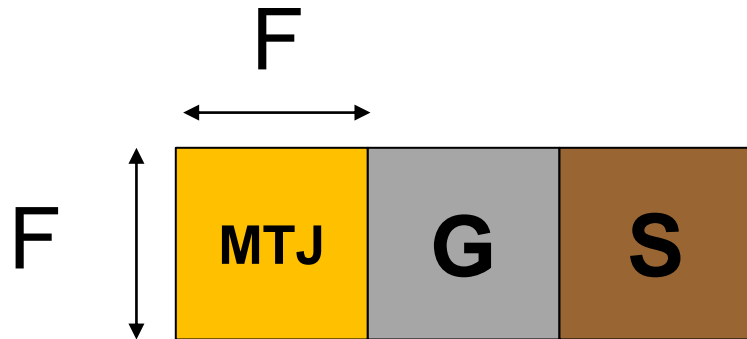
Write current

- Critical current for switching magnetization: $J_C = \left(\frac{\alpha}{\eta}\right) \left(\frac{2e}{\hbar}\right) M_s H_k t + 2\pi M_s$
 - α = damping factor
 - H_k = anisotropy field
 - t = free layer film thickness
 - M_s = saturation magnetization
 - η = STT efficiency parameter
- Current costs energy!
→ Want as low as possible



2 min Exercise: Current supply

- Typical value $J_C = 3 \text{ MA/cm}^2$
- 10nm node Intel transistor: $F = 50 \text{ nm}$. $I_{on} = 1 \text{ mA}/\mu\text{m}$ @ $W = 1F \rightarrow I_{on} = 50 \mu\text{A}$.
- **For a $1F^2$ MTJ, can this transistor supply sufficient current density?**



2 min exercise – The write current

1. What parameters should we use to decrease the write current?
2. How can we practically change these?

$$J_C = \left(\frac{\alpha}{\eta}\right) \left(\frac{2e}{\hbar}\right) M_s H_k t + 2\pi M_s$$

η = STT efficiency parameter

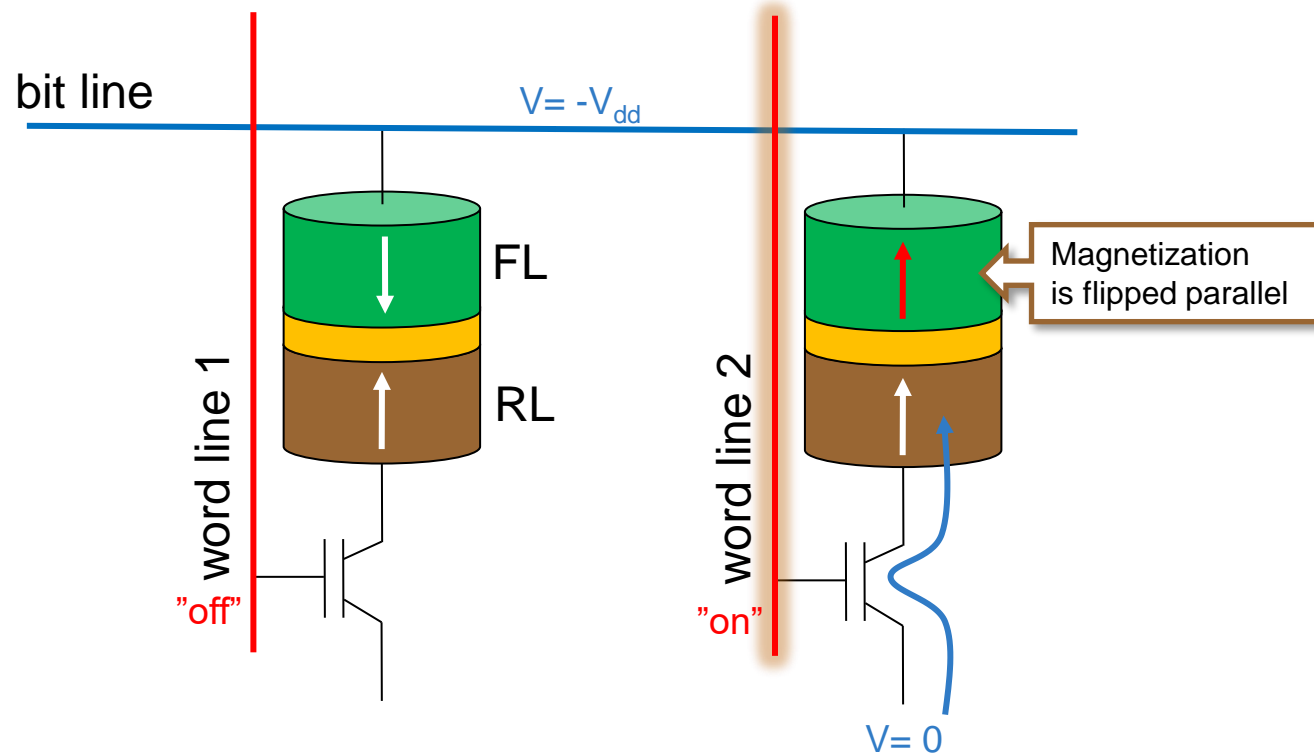
α = damping factor

M_s = saturation magnetization

t = free layer film thickness

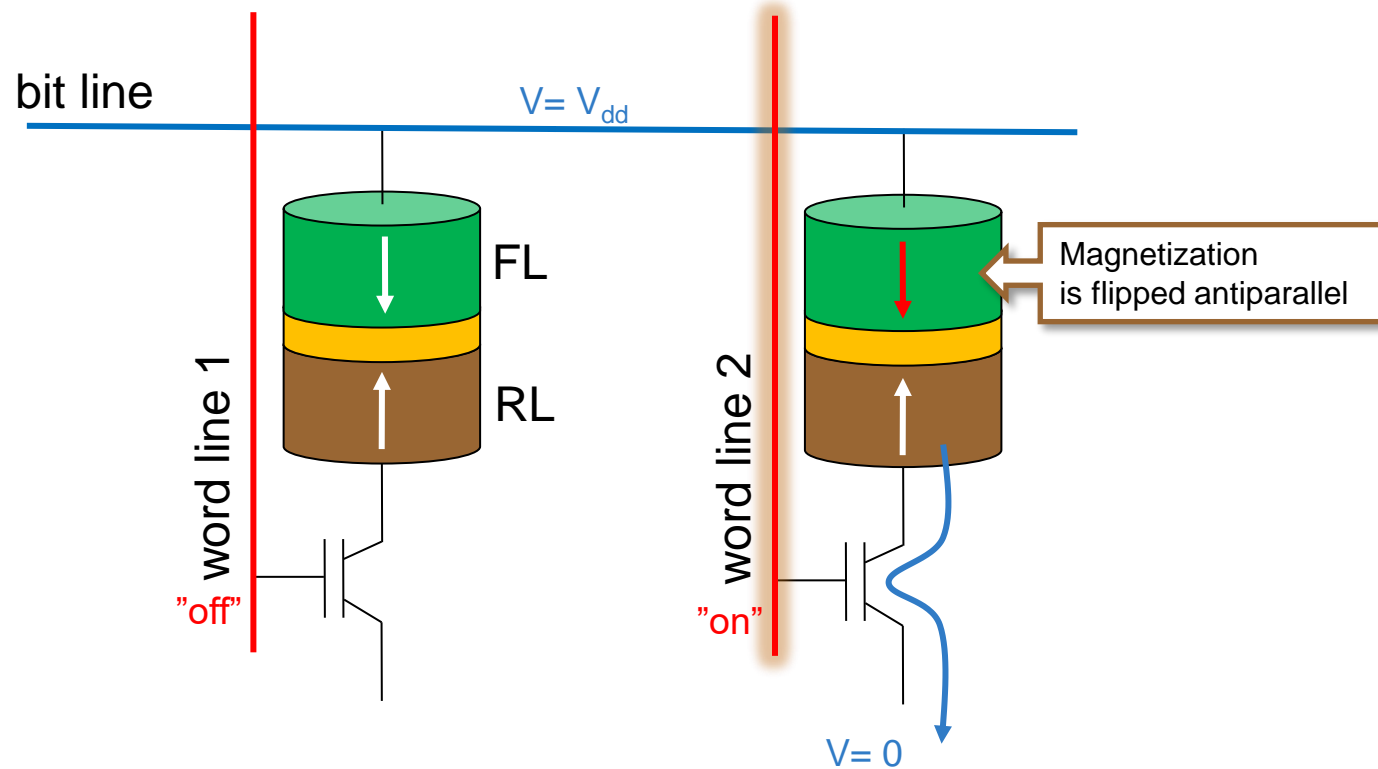
H_k = anisotropy field

STT-MRAM Write "0"



- "0" is low resistive state (parallel)
- Word line chooses device
- Bit line biased negative \rightarrow current "upwards" \rightarrow parallel spin

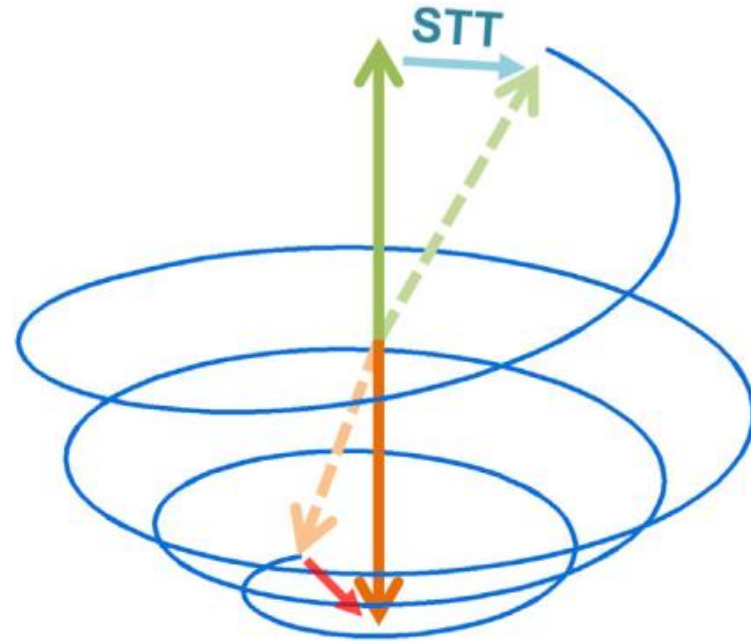
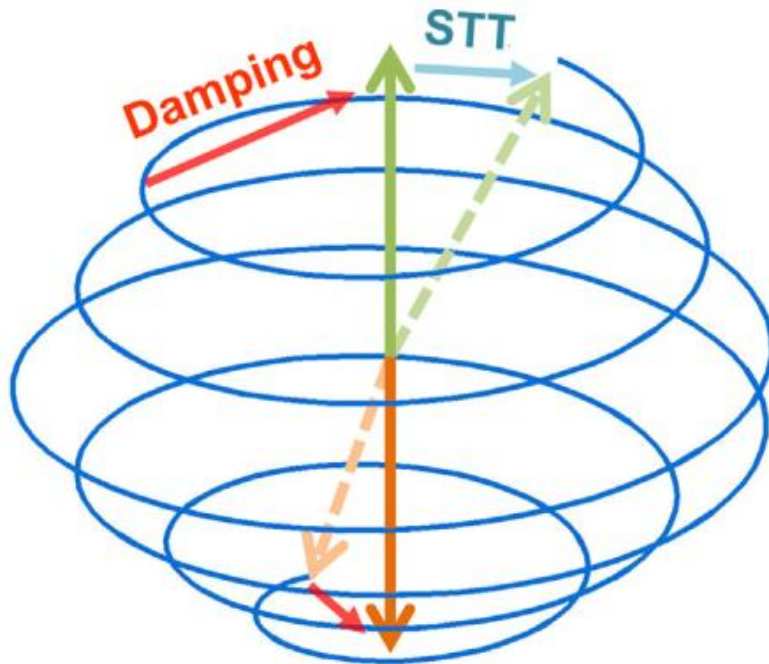
STT-MRAM Write "1"



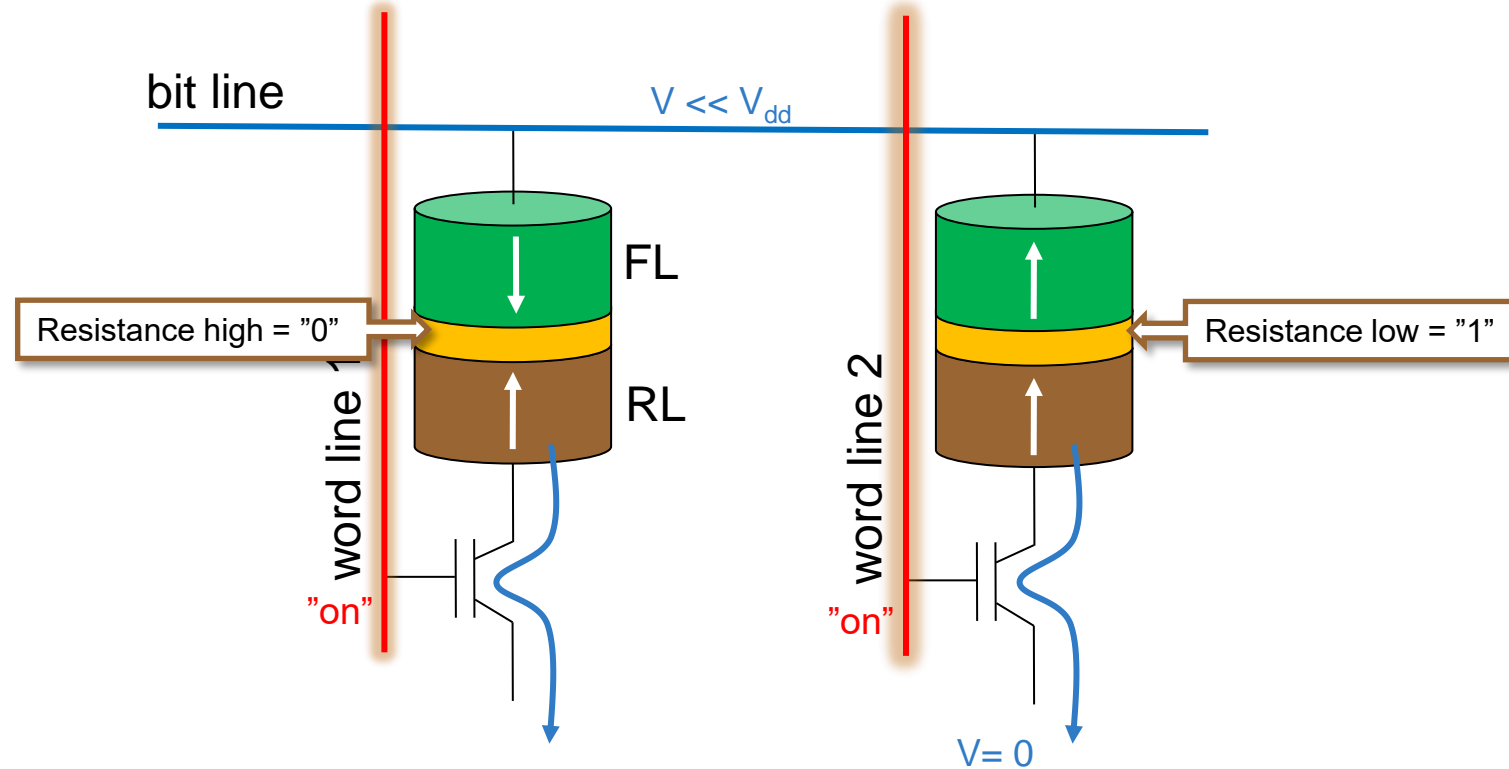
- "0" is low resistive state (parallel)
- Word line chooses device
- Bit line biased negative \rightarrow current "upwards" \rightarrow parallel spin

Damping: Write current and speed

- Larger damping factor α gives faster switching due to faster to achieve end magnetization.
 - But also means larger switching current is needed!
- α gets stronger with spin-orbit coupling ($\propto Z$).
- α gets stronger with imperfections (roughness, defects, interfaces, surrounding layers)



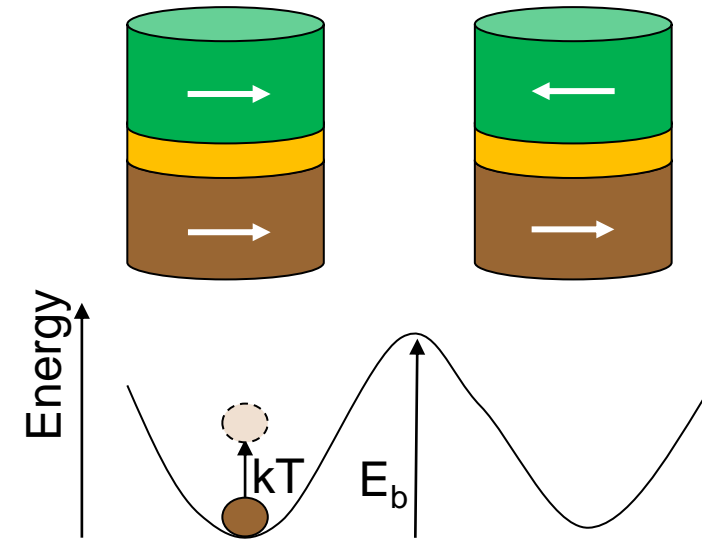
STT-MRAM Read



- Read current much lower than needed to flip magnetization
- Parallel / Antiparallel spin in free layer \rightarrow low / high resistance i.e. data read out

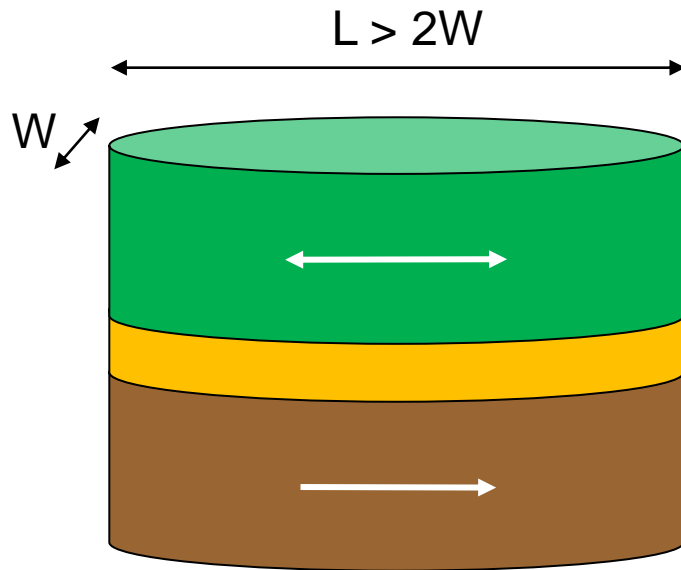
Data retention and scaling

- Free layer needs to retain its magnetization despite disturbance
- Typical target is 10 years retention $\rightarrow 60 \text{ kT} \sim 1.5 \text{ eV}$ energy barrier for switch
- Possible disturbance:
 - Thermal fluctuations
 - Read event
 - Typical prob. for unwanted switch $\sim 10^{-21}$
- $E_b \propto K_u V \rightarrow$ scaling down decreases barrier height!
 - Need very high K_u for ultra-scaled memory cells
 - This was for long a major road block

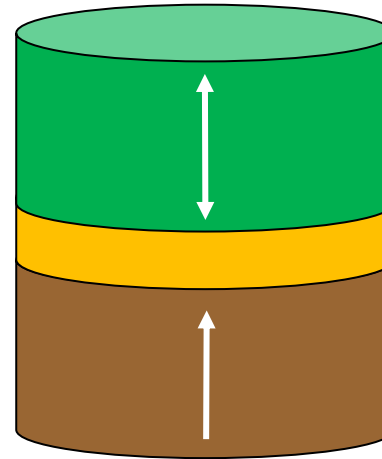


In-plane or perpendicular magnetization?

- Magnetization can also be *out-of-plane* (perpendicular)
- In-plane relies on *shape anisotropy* → not scalable beyond 60 nm
- Perpendicular Magnetic Anisotropy does not rely on shape anisotropy

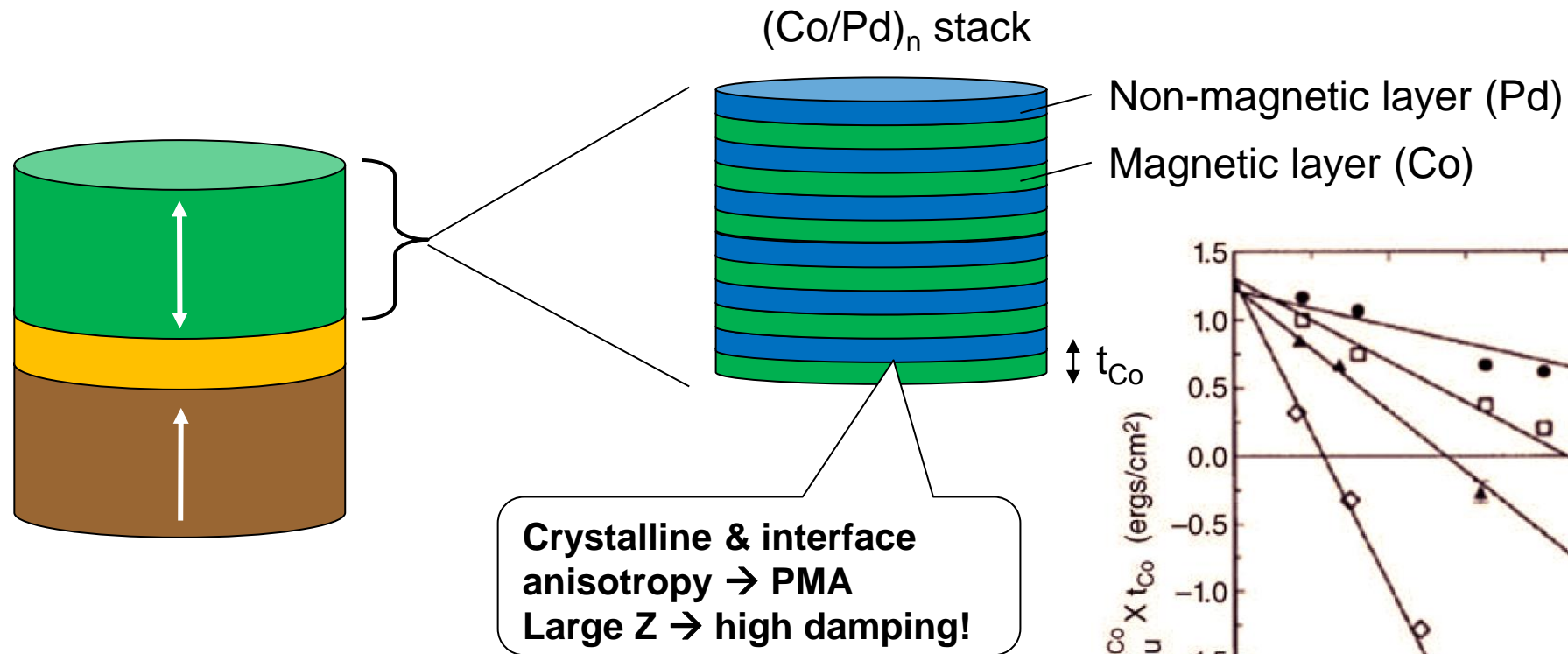


In-plane

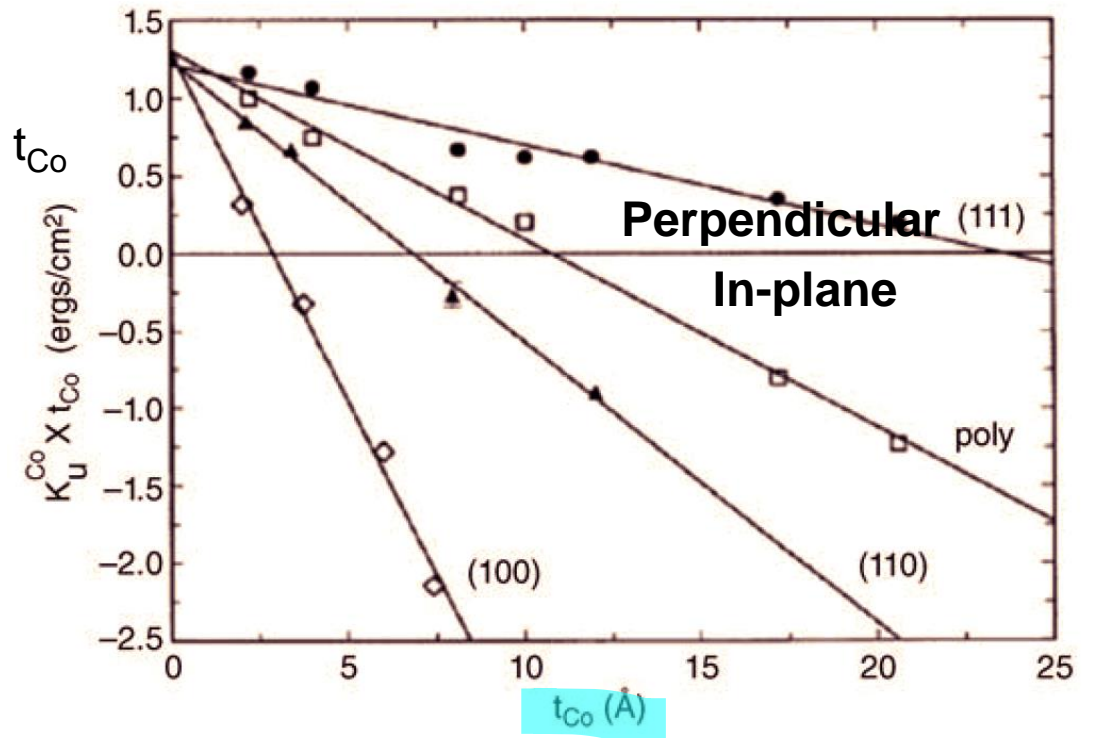


Perpendicular
(to the) Plane

Designing for Perpendicular anisotropy



PMA also observed in very thin CoFeB layers
 → PMA with less damping! $J_c \sim 3 \text{ MA/cm}^2$

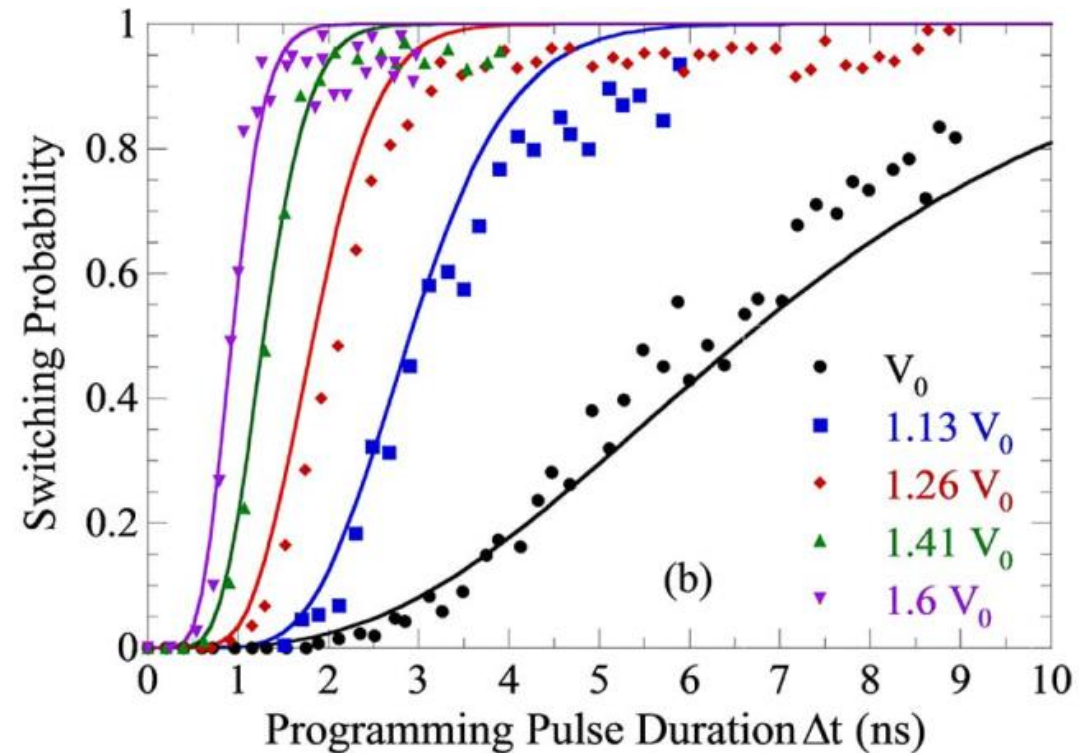
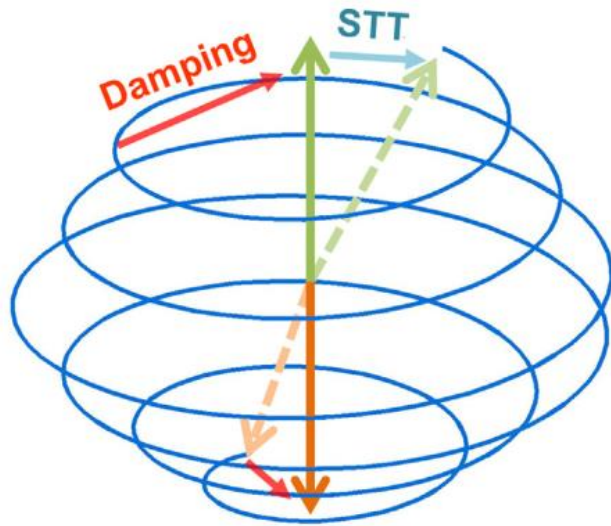


STT-MRAM as storage

	DRAM	3DNAND	RRAM	PCM	STT-MRAM
Nonvolatile	No	Yes	Yes	Yes	Yes
Speed (ns)	10	10^4	< 5 ns	10 ns	< 5 ns
Energy use (pJ/write)	0.1	1	0.1-1	>1	< 0.2 pJ
Endurance (cycles)	10^{16}	10^5	10^6 - 10^7	10^9	$>10^{15}$
Multilevel?	No	Yes	3-6 bit	Yes	No
Scalability	6-8F ²	3D!	3D!	3D!	6F ²
Other	Destructive Read	High Voltage	Abrupt SET	R drifts	Scaling limited by needed current

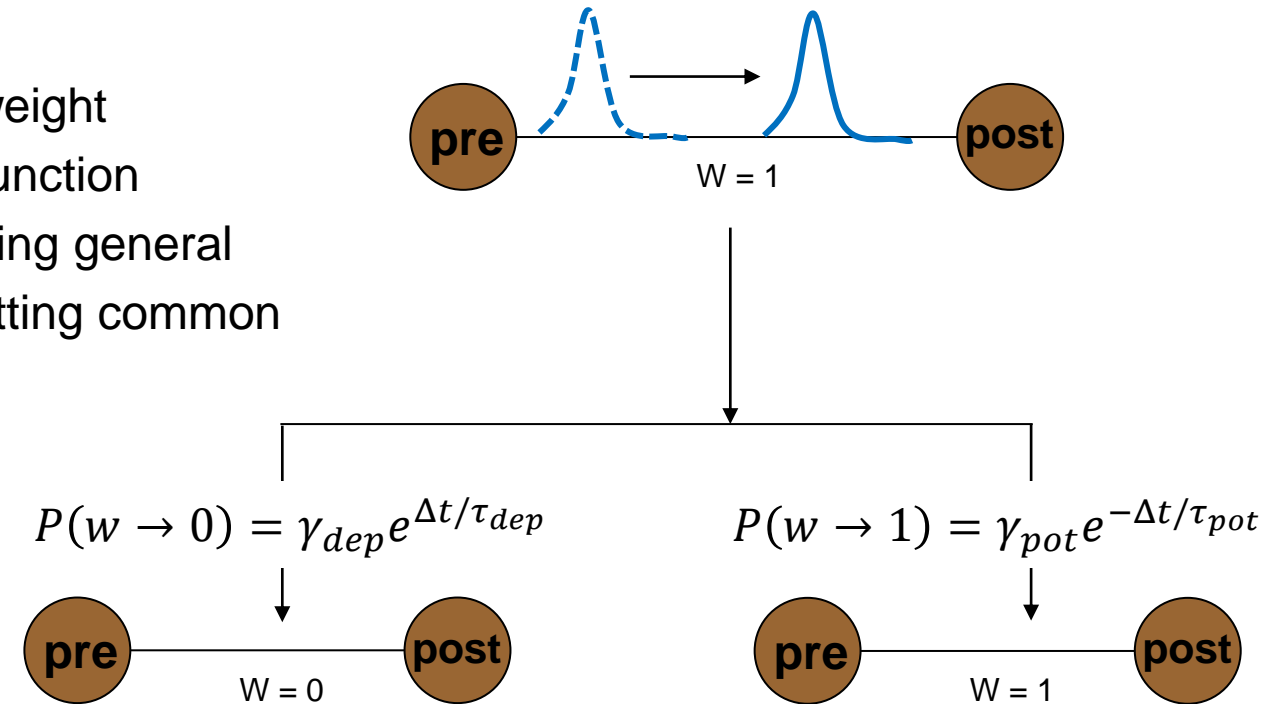
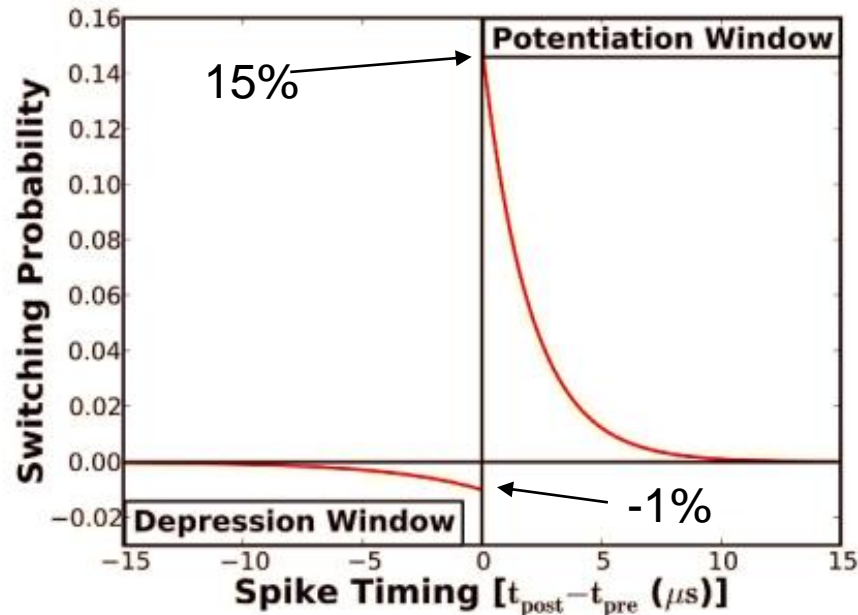
MRAM in SNNs

- Binary device, no memristor
- With J_c near threshold \rightarrow M switching becomes stochastic!



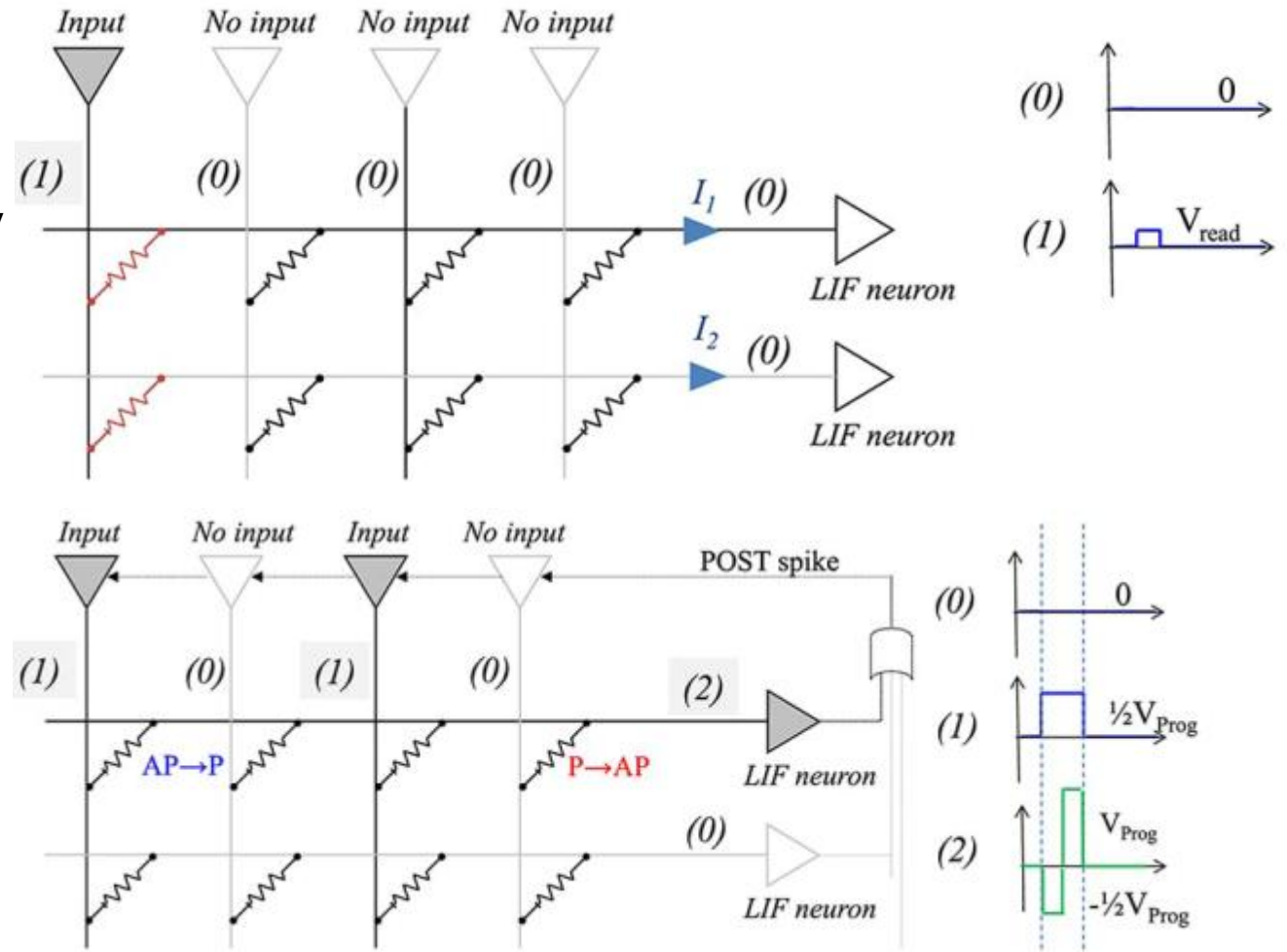
Stoichastic STDP

- Stoichastic STDP
 - Every post-spike can lead to flip of synapse weight
 - Switching probability P_{switch} given by STDP function
 - Potentiation prob. not too high \rightarrow keeps learning general
 - Depression prob. much lower \rightarrow avoids forgetting common features between classes



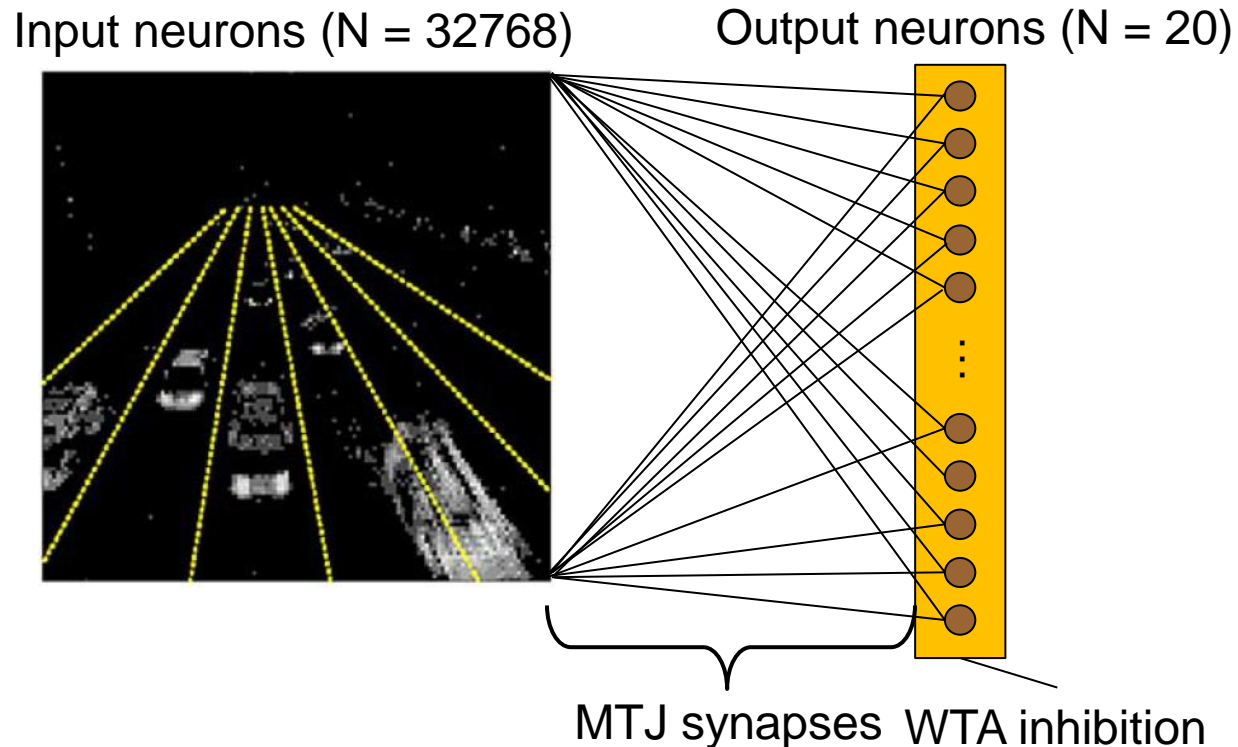
Example of implementation

- 1 layer network.
- CMOS LIF neurons on output.
- Integrates signals from parallel M_s devices only
→ Digital weight representation
- Upon post spike:
 - Post-neuron applies (2)
 - All recently active pre-neurons apply (1)
 - Inactive preneurons apply (0)
 - Lateral inhibition of output neurons (WTA)
- (1) + (2) give certain probability to flip M_s
 - Set to equal to 10% in the paper.

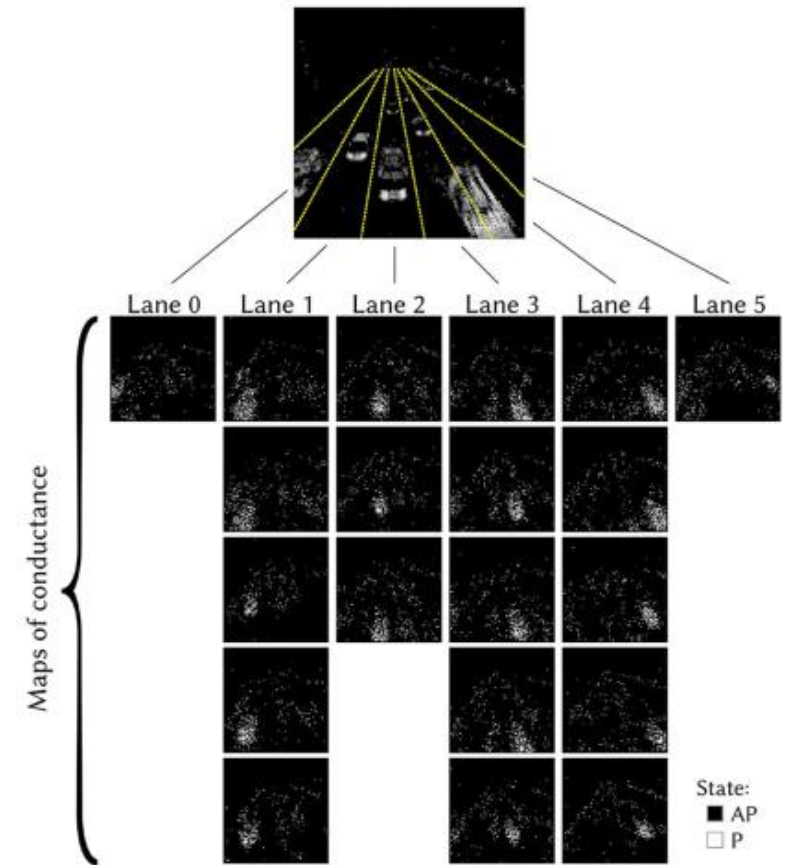


Test of implementation (car counter)

- Pasadena 5 lane highway recorded with neuromorphic retina
 - No frames! Spikes when intensity of pixel changes



Natural specialization on particular lanes!



Detection accuracy in lanes 1-4 ~ 97%
Power consumption for learning 180 nW!

Effect of device variability

- Variations in min and max resistance \rightarrow variation in current \rightarrow dramatic variation in switching probability!
- Synaptic variations (SV) of 5, 10 and 25% (very high!)
- Very robust to device to device variations!

