

## ES14 - ADVANCED ADCs

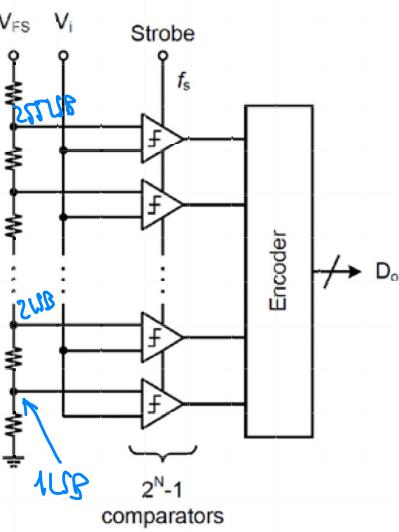
30/11/2021

- Goals: - reducing  $T_{conv}$  (increase speed) } serial (pipeline) processing  
 - reducing # of components } parallel processing  
 ↳ which means reducing Area & PC

### SUBRANGING ADC

BFM: In the Flash ADC we create "all" possible voltage values from  $V_{ref}$  to  $V_{FS}$  using resistors ( $2^n$ )

Then, using ( $2^n - 1$ ) comparators, we compare  $Ain$  to all these just created values



From a certain point downwards the comparators will be triggered b/c  $Vin$  is higher than their respective  $Vref$  while the previous comparators will not be b/c  $Vin$  is lower than their respective  $Vref$ . Then, changing  $Vin$ , some comparators will start triggering and others will stop triggering.

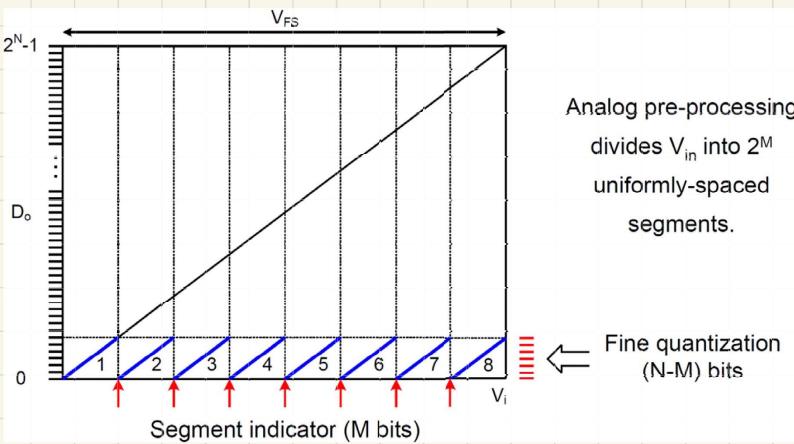
Let's suppose to have  $V_m = 2.5V = FSR/2 = Vref/2$  so we'll have 122 comparators that are already triggered and 122 comparators that are not triggered yet.

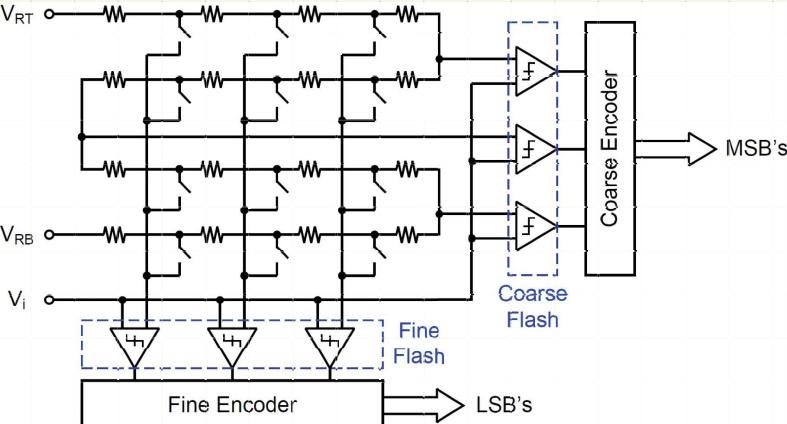
The comparators that matter are those ones that are very close to  $Vm$ , b/c they are those ones on which the encoder relies to decide the output code  $D_o$ .

↓  
 The useless 122 comparators that are already triggered continue to dissipate power → ISSUE: POWER CONSUMPTION

### SOLUTION: SEGMENTED QUANTIZATION

basic idea: Let's use a lower # of comparators (i.e. 16 instead of 256) in such a way that they will give us a coarse quantization so they will give us just the M bits of our conversion and then once we decide where we are in the conv., we'll turn on the other 16 comparators, in such a way in each sub-sample we have 16 LSBs, so eventually our 256 LSBs can be split in 16 main-levels and each main-level will decide which subset of comparators have to be turned on in order to complete the fine quantization of our analog input voltage  $Ain$ .





At the beginning the switches are all open and only the coarse flash comparators are connected to the ladder.

Once the coarse comparators provide the coarse conversion, so once the coarse encoder understands where the  $V_{in}$  approximately is (MSB conversion).

Only after the MSB conversion the coarse encoder will sequentially enable the switches: one is closed, while the others remain open. As a consequence, the fine comparators will turn on providing the fine conversion: the fine encoder understands where  $V_{in}$  exactly is!!!

We subdivided our FSR in, say, 16 levels and once we provide those 4 MSBs we understand in which of those 15 levels we are and so we turn on the remaining 16 switches in that specific macro-block (the macro-block depends on the result of the MSB conversion).

ADVANTAGE: we do not need 256 comparators but for an 8 bit ADC we need 4 bits for the MSB (coarse) conversion and 4 bits for the LSB (fine) conversion.

→ So in total we need 16 comparators for the MSB conv. (actually 15 bcz one is always connected to the lowest possible value) and 16 comparators for to convert the LSBs

⇒ we only need 32 comparators instead of 256!! 😊

What does it mean?

- Area drastically reduced!
- Power consumption drastically reduced!

REMAINING ISSUES: we still need 256 resistors

MAYBE?? NOT! {

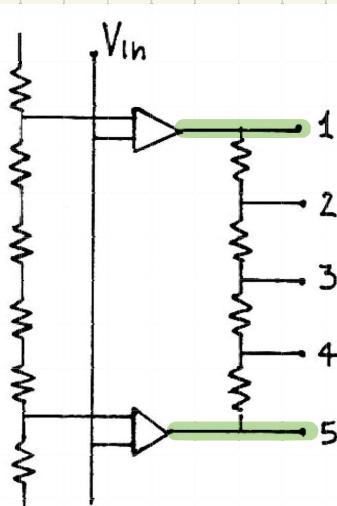
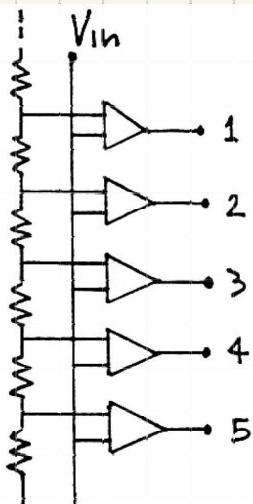
- Compared to a standard flash ADC, this subranging ADC takes almost 3 times to perform the conversion

↳  $T_{conv, subranging} \approx 3 T_{conv, flash}$

The 16 fine comparators should have  $V_{os, fine} \ll 1LSB = 19\text{mV}$  (it should be  $V_{os, fine} \leq 1\text{mV}$ )

## INTERPOLATING FLASH ADC

Let's remove 3 comparators from every 4 comparators:



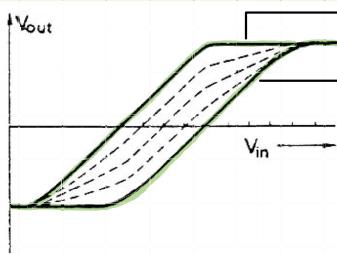
This means that for a 8 bit ADC we'll have 64 comparators instead of 256

Only 64 lines will reach the encoder that will produce a 6 bit ( $2^6 = 64$ ) code

Now, let's add resistor b/w our comparator output and the following comparator output

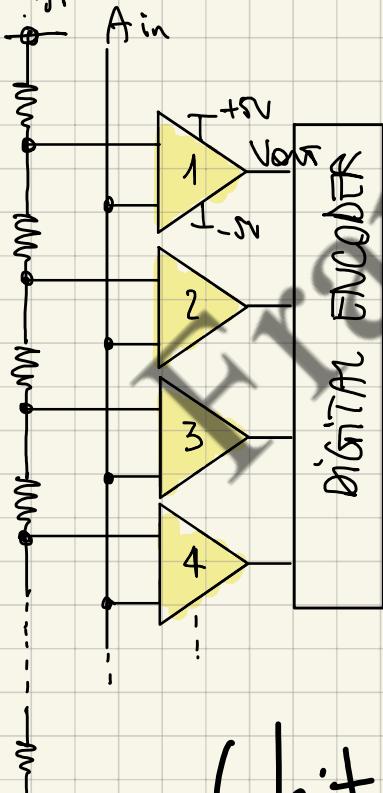
These added resistors will mimic the comparators we have removed

→ The 3 missing outputs are reconstructed w/ a network of resistor b/w the remaining outputs



$V_{out}$  is an index of similarity b/w  $V_{in}$  and  $V_{ref}$

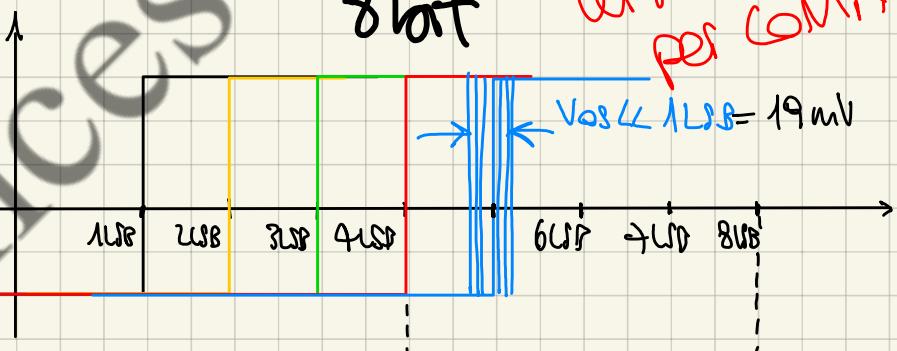
In a typical Flash ADC we have:



8 bit

let's 25 MOSFETS per COMPARATOR

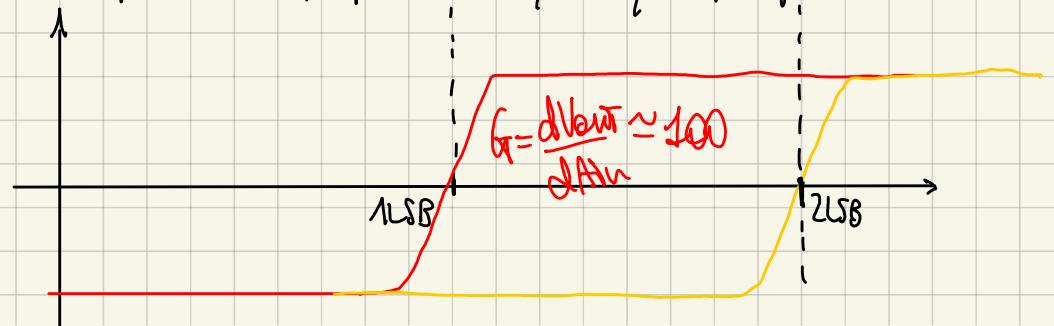
$$V_{os} \ll 1LSB = 19mV$$



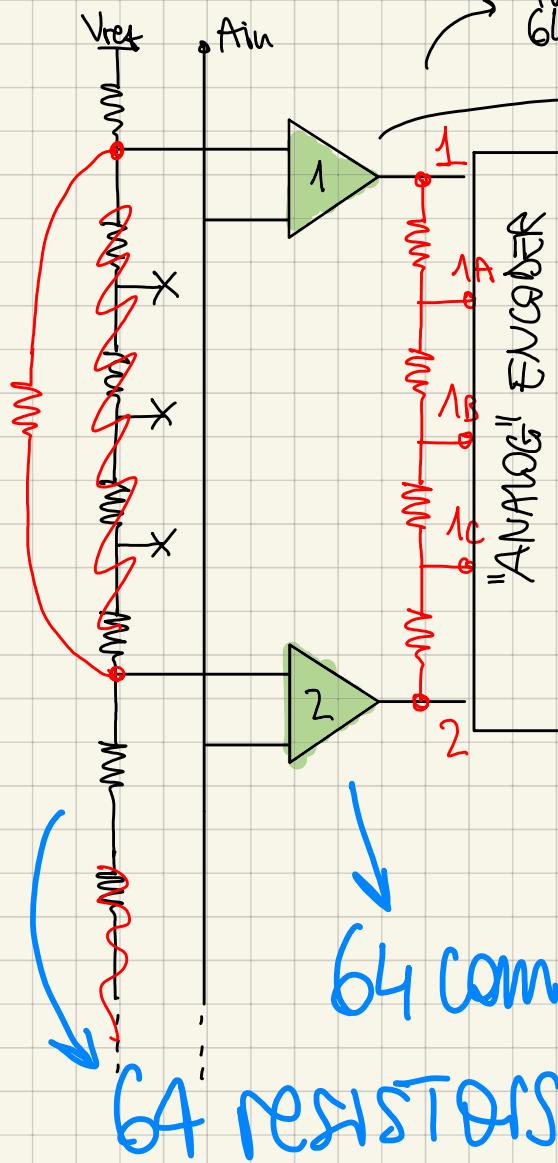
$$\frac{dV_{out}}{dA_{in}} = \frac{10V}{(2^8 \cdot 1LSB)} \approx \frac{10V}{1mV} \approx 10k \rightarrow \text{very high Gcomparator}$$

This means that in a Flash ADC the higher the # of bits and so the # of comparators, the higher the # of resistor we must put inside each comparator to get a very very high gain

6 bit

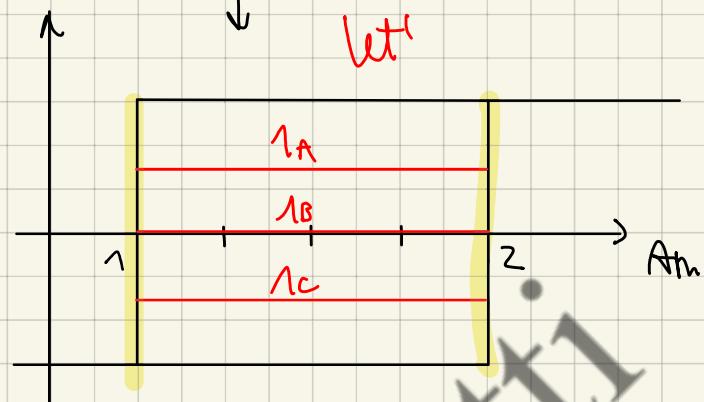


Now the idea is:

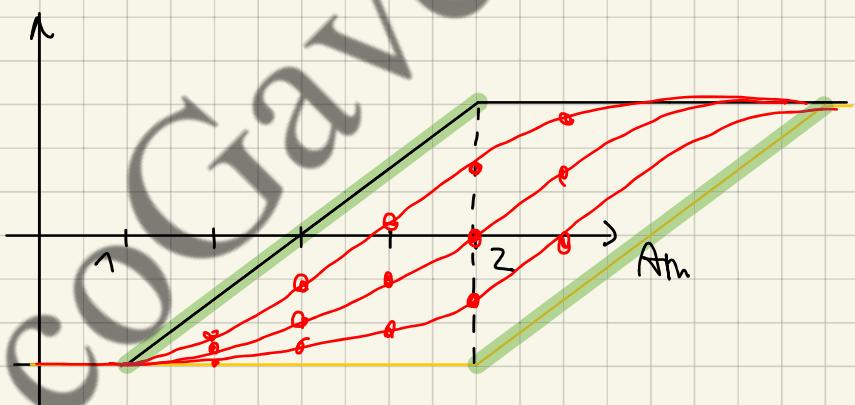


This behaves as a 6 bit ADC once it implements 64 comparators

if we use the same comparators as before we end up w/ an error



Let's use much finer comparators:



$$G = \frac{dV_{out}}{dA_{in}} = \frac{10V}{1LSB_{6bitADC}} = \frac{10V}{\frac{5V}{2^6}} = \frac{10V}{\frac{5V}{64}} = \frac{1}{4} G_{6bit}$$

Now we have the same # of outputs and what we have to do is simply check if one output commutes or not.

In a flash ADC, when one comparator commutes, the output was either  $+5V$  or  $0V$ .

Here instead we relax the pain and we create analog intermediate voltages, so now the output is an analog output

The output of a flash ADC should feed a standard digital encoder (where the components of the encoder are AND gate, OR gate and so on), here instead, the outputs should feed an ANALOG ENCODER which checks if the voltage is above or below  $\Delta$ , so checks which comparators are triggered and which are not, also considering the intermediate analog voltages

⇒ The "ANALOG" ENCODER checks the zero-crossing!!!

### ADVANTAGES:

①  $\frac{1}{4}$  of comparators wrt a standard flash ADC are used

↳ at least  $\frac{1}{4}$  less of power consumption

② Since we don't need to have a very high gain anymore, we don't need 25 transistors per comparator. Let's say we need the gain to be very soft, so each comparator will have a lower # of transistors.

↳ much lower amount of power gets dissipated

→ also the input parasitic capacitance will be much more reduced b/c the offset can be more relaxed

③ HUGEST AMOUNT OF AREA SAVING!!!

④ The dynamic performances are improved

In a flash ADC when we vary Ain we have to wait until the input of each of the 256 comparators settles to the new Ain value.

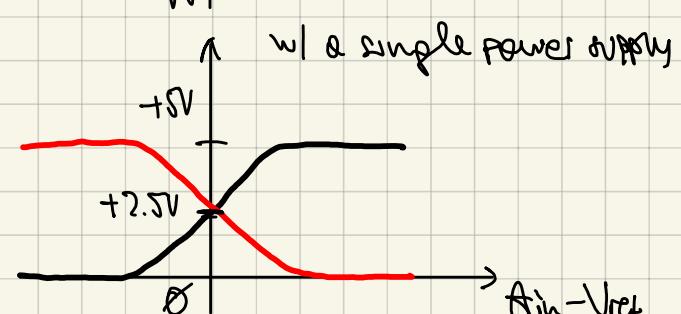
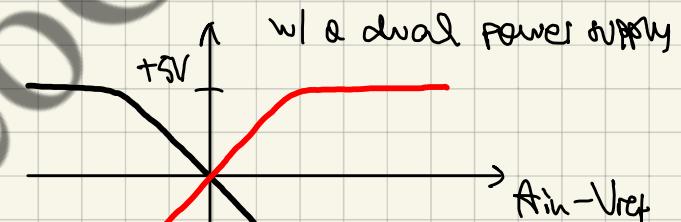
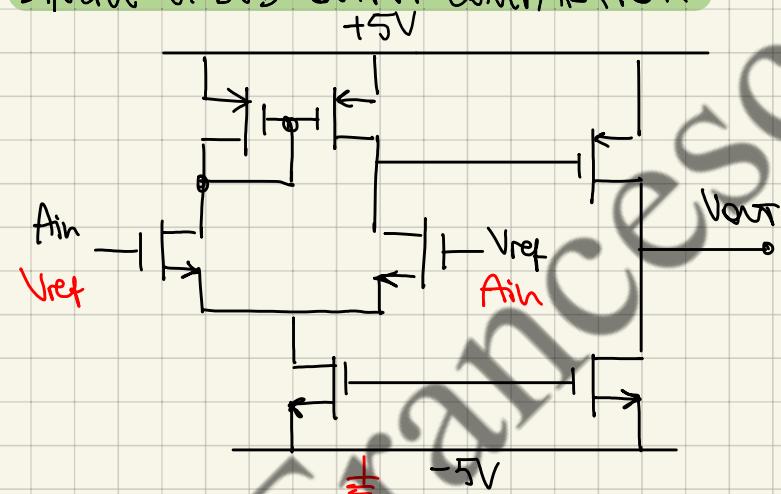
The interpolating ADC could work faster than a flash ADC in principle, but instead don't b/c the gain gets relaxed, we use a much more complex ENCODER.

A further improvement could be using a DIFFERENTIAL COMPARATOR and two different loaders at the output.

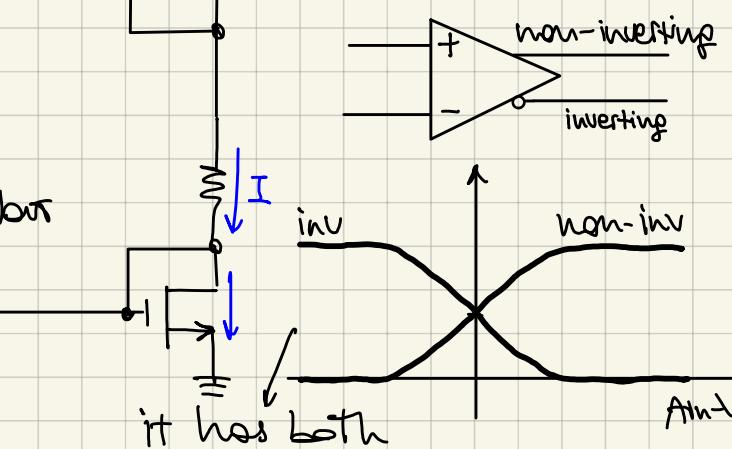
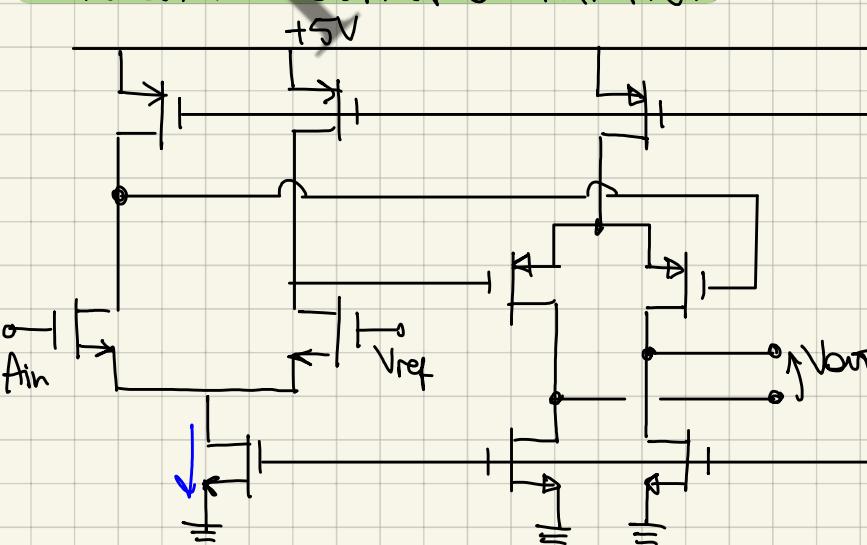
Let's proceed step by step.

Up to now we considered a comparator like this:

SINGLE ENDED OUTPUT COMPARATOR:



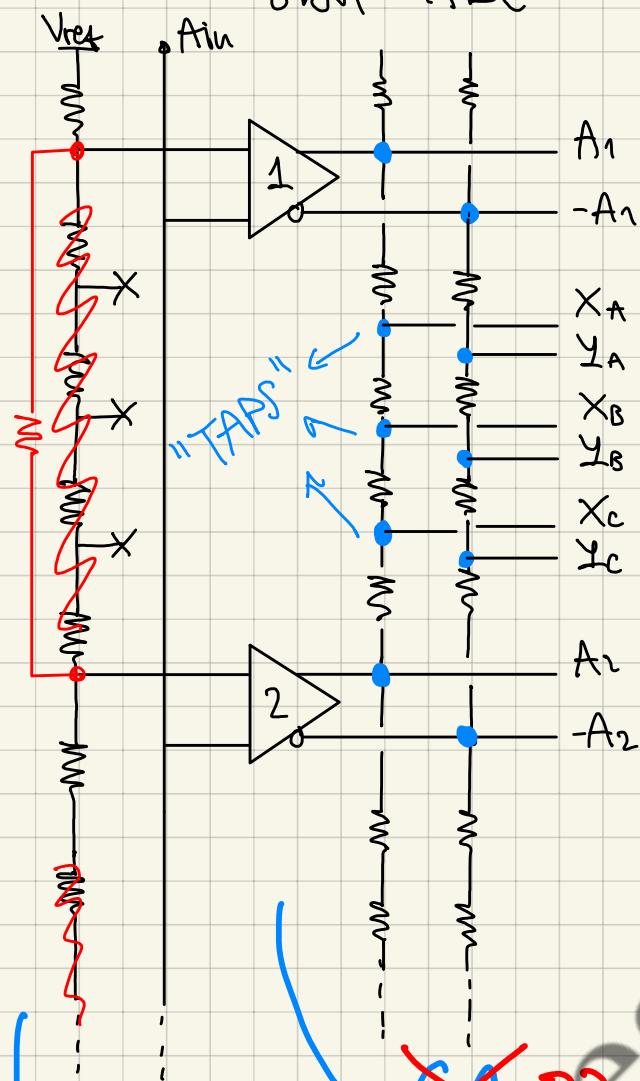
DIFFERENTIAL OUTPUT COMPARATOR



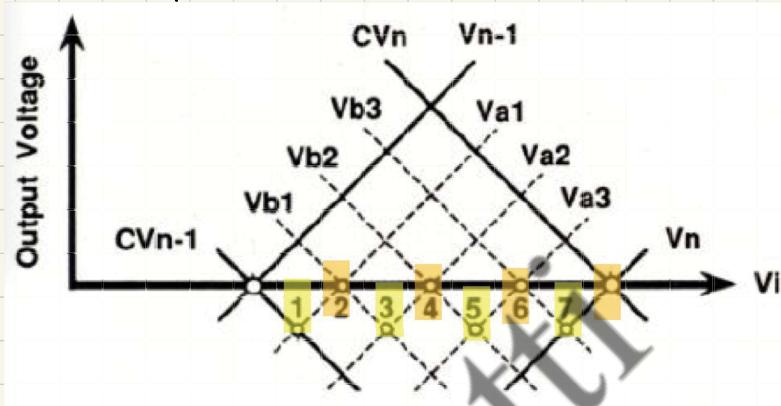
it has both

ADVANTAGE: Using differential comparators we can further halve the number of comparators and so also of the input resistors

### 8-bit - ADC



→ here we have 2 sets of 4 interpolation resistors



Interpolation resistors produce intermediate levels: 4 resistors → 2 intermediate levels

w/ 4 resistors per side, we have:

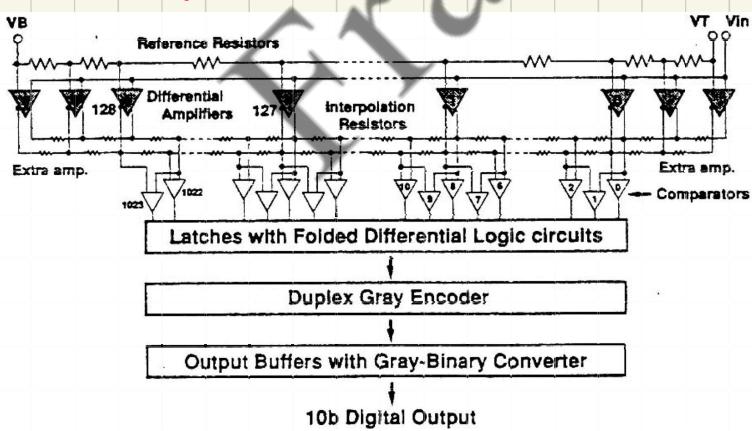
$M=8$  INTERPOLATION FACTOR

↳ corresponds to 3 EXTRA BITS

~~64 DIFFERENTIAL "COMP.s"~~

~~64 INPUT RESISTORS~~

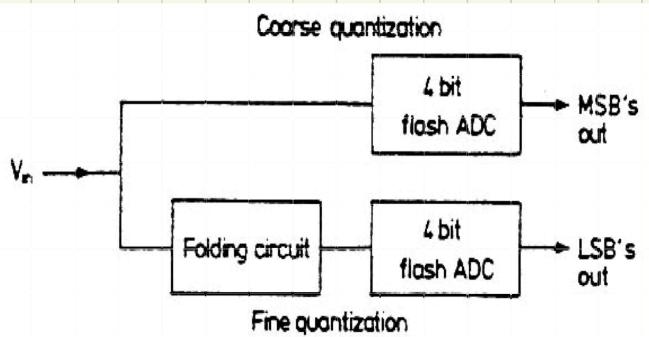
EXAMPLE: INTERPOLATING 10-ADC w/ DIFFERENTIAL COMPARATORS



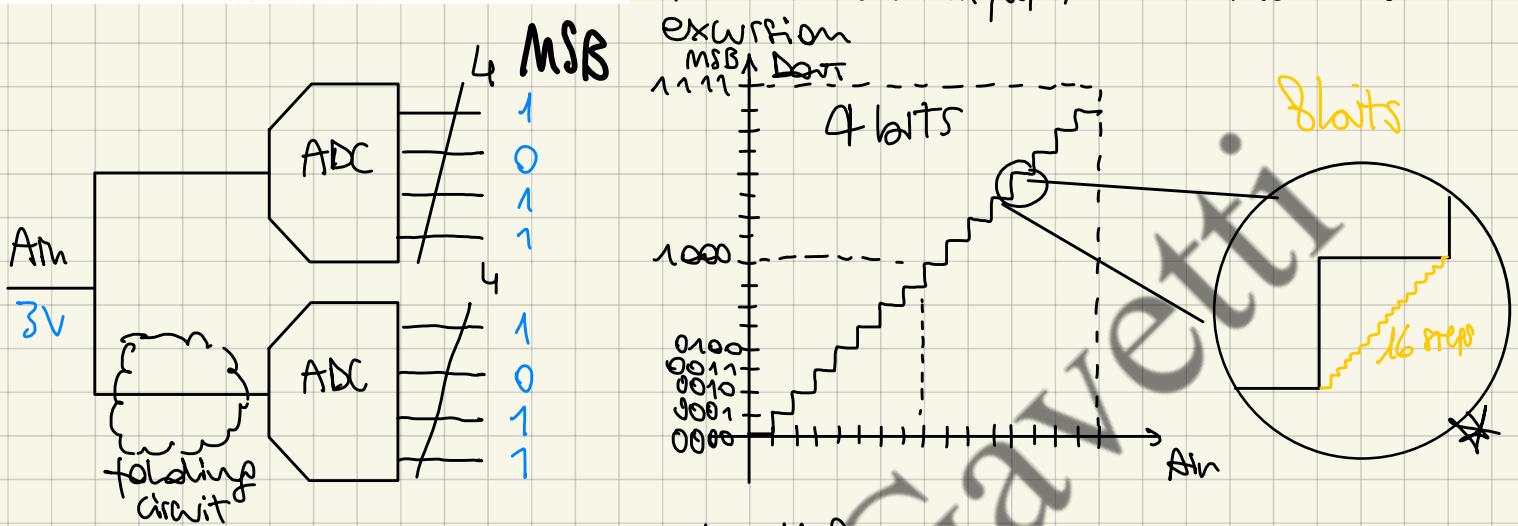
$2^{10} = 1024 \Rightarrow$  using a standard flash ADC we should use 1024 input resistors and 1024 standard comparators

Instead, using an interpolating ADC w/ differential output comparators we'll use  $1024/8 = 128$  input resistors and 128 comparators

# FOLDING ADC (FLASH FOLDING)

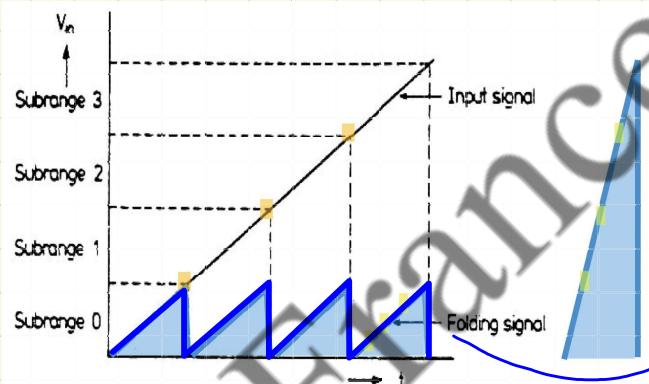


**BASIC IDEA:** by means of an analog signal pre-processing we can separately obtain the MSBs and the LSBs. The former are obtained w/ a low resolution ADC whereas the latter is obtained by folding the input signal into the folding circuit that translates the input ramp into a saw tooth, w/ limited maximum excursion



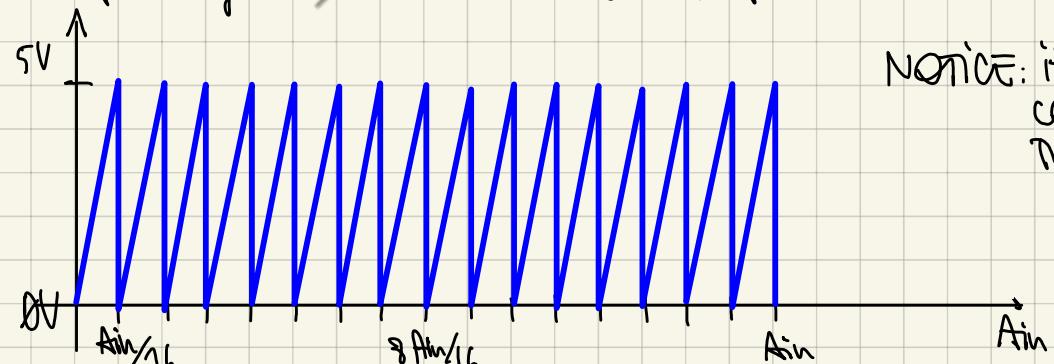
\* if we want to subdivide each step by 16 levels we need something that converts Ain into the error b/w the coarse quantization and the straight line

What we need is a **FOLDING CIRCUIT** which is a circuit whose input-output relationship is the blue one:



If the folding circuit is able to split the range b/w 0 and 5V into 16 folds, then the height of the folding circuit will be the error b/w Vin and the quantization achieved w/ the 16 levels ADC

The folding circuit should be able to provide this in-out conversion

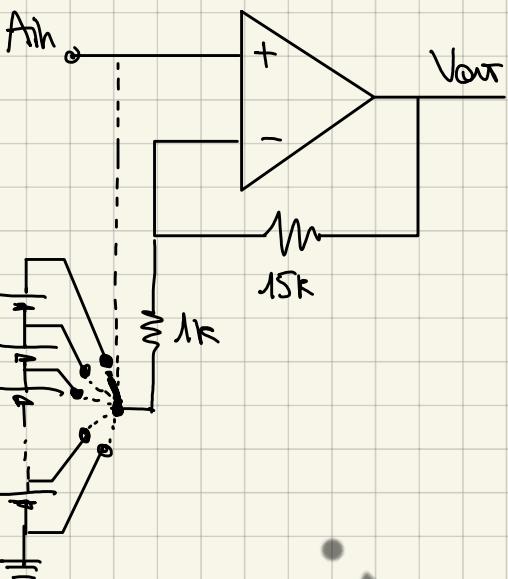


NOTICE: it gives us the error b/w a coarse quantization and the actual Ain

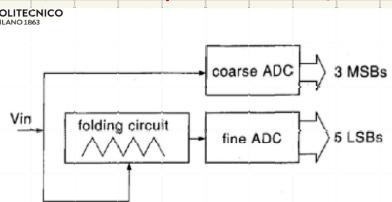
Notice: This solution requires  $2^h - 1 + 2^h - 1 = 15 + 15 = 30$  comparators

How could we design the folding circuit?

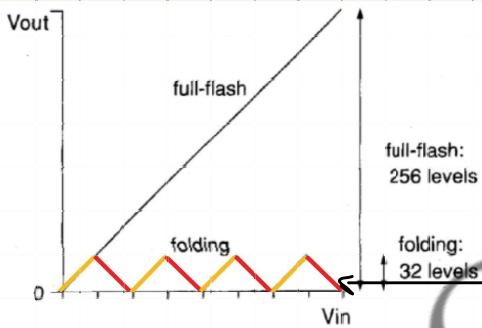
In principle, it's easy, it could be:



## ANOTHER FOLDING FLASH ADC

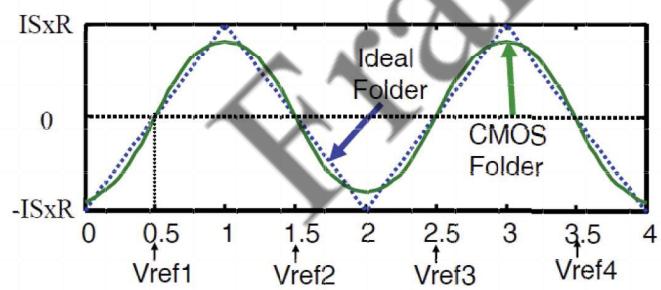
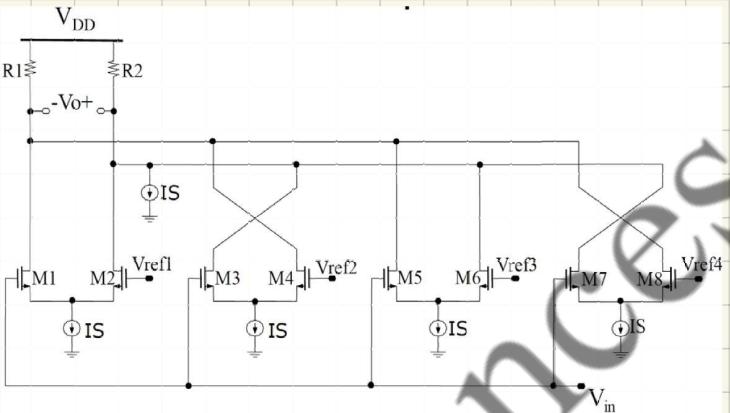


$$2^3 \cdot 1 + 2^5 \cdot 1 = 38 \text{ comparators instead of } 256$$



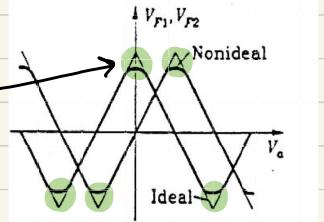
The folding circuit instead of being a kind of sawtooth, is

## Possible FOLDING FLASH ADC IMPLEMENTATION:



Instead of using just one folding circuit, we could use two of them: one shifted by the other by  $1/2$  LSB

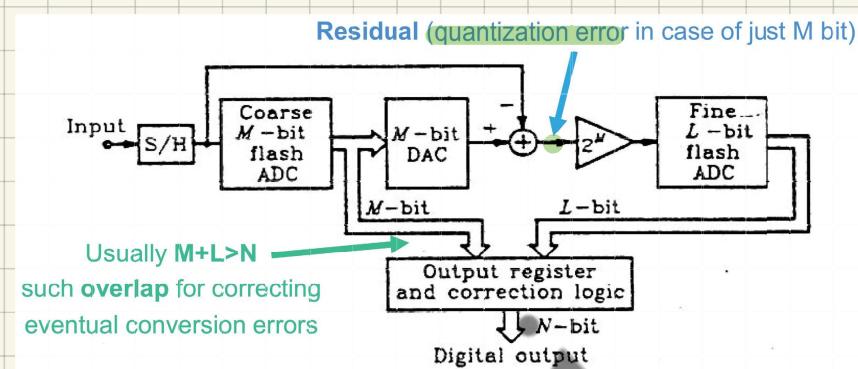
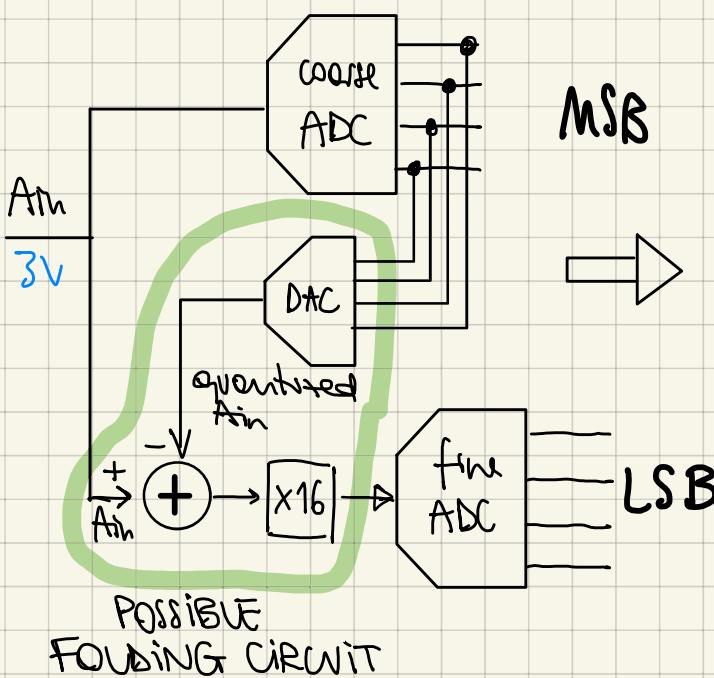
↓  
in this way we solve this issue  
of the smoothing in the curve  
but we'll use the LSB coming from  
only one of the two  
which one? The one that  
gives the lower LSB



## HALF-FLASH ADC

it comes from the folding ADC

We've seen that the folding ADC provides us the error of the coarse quantization



During the 1st step the signal is applied to the 1st flash converter which makes a low resolution conversion (worse quantization) obtaining the M-MSBs.

These are converted by a DAC, whose accuracy must be at least equal to that one of the whole ADC, ie N bits.

This value, which is an approx. of the input signal, is subtracted from the same input signal and the result is the RESIDUAL, ie the error committed approximating the signal w/ M bits.

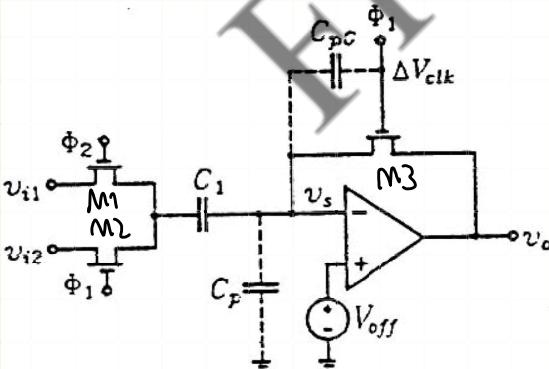
The residual is amplified and sent to the 2nd flash ADC that determines the L-LSBs doing the FINE QUANTIZATION.

The outputs of the converters are summed and the result is the N bits digital output.

We desire to have  $M+L>N$  in order to perform the correction of any conversion error  $\rightarrow$  The additional bits are named OVERLAP BITS

## INTERNAL DESIGNS OF HALF-FLASH ADC

### COMMUTATING AUTO-ZEROING TECHNIQUE



1st phase: { M1, M2 open  $\rightarrow$  The comparator acts as an Opamp connected as a buffer  
M3 close

$\hookrightarrow$  Then thanks to the FB,  $v_{os}$  is applied to  $C_1$

2nd phase:

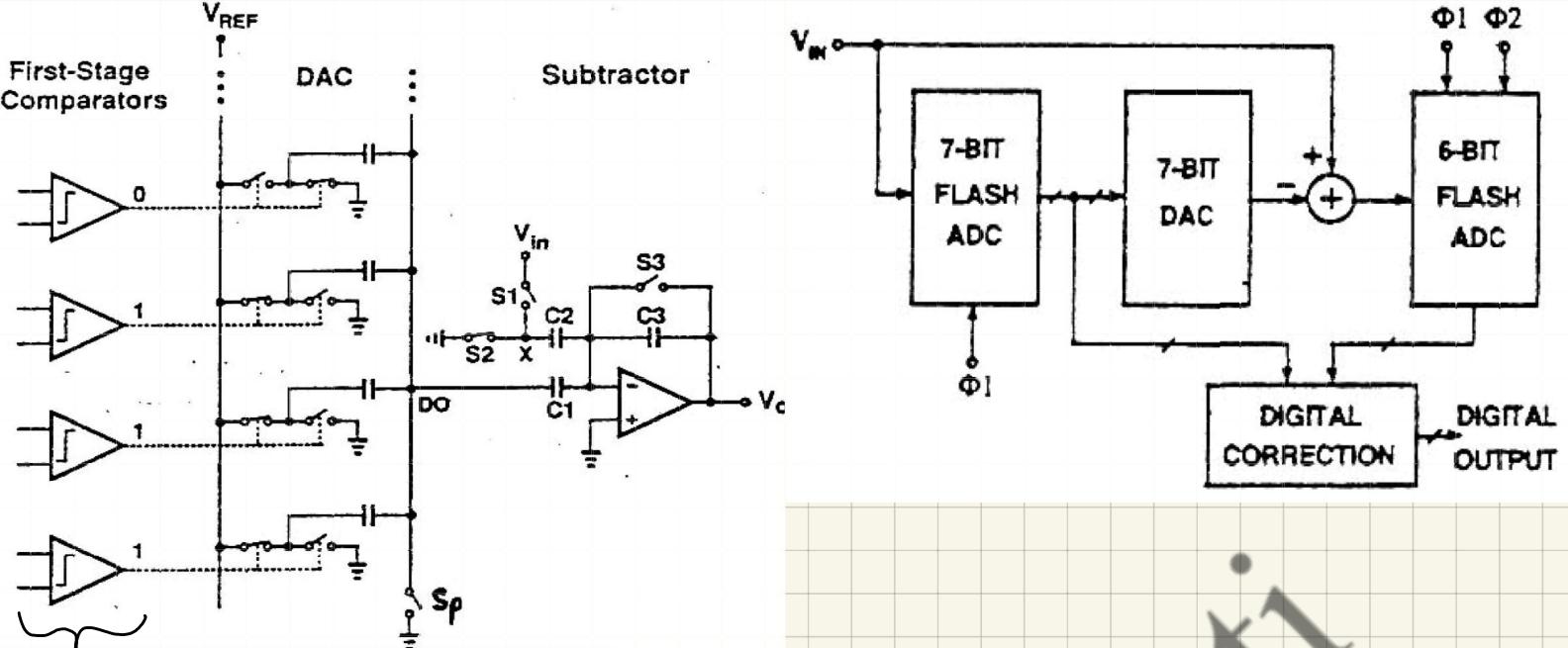
So closing M2 and applying  $V_{il} = 0$ , we store  $v_{os}$  across  $C_1 \rightarrow$  AUTO-ZEROING

3rd phase { M1 close  
M2, M3 open

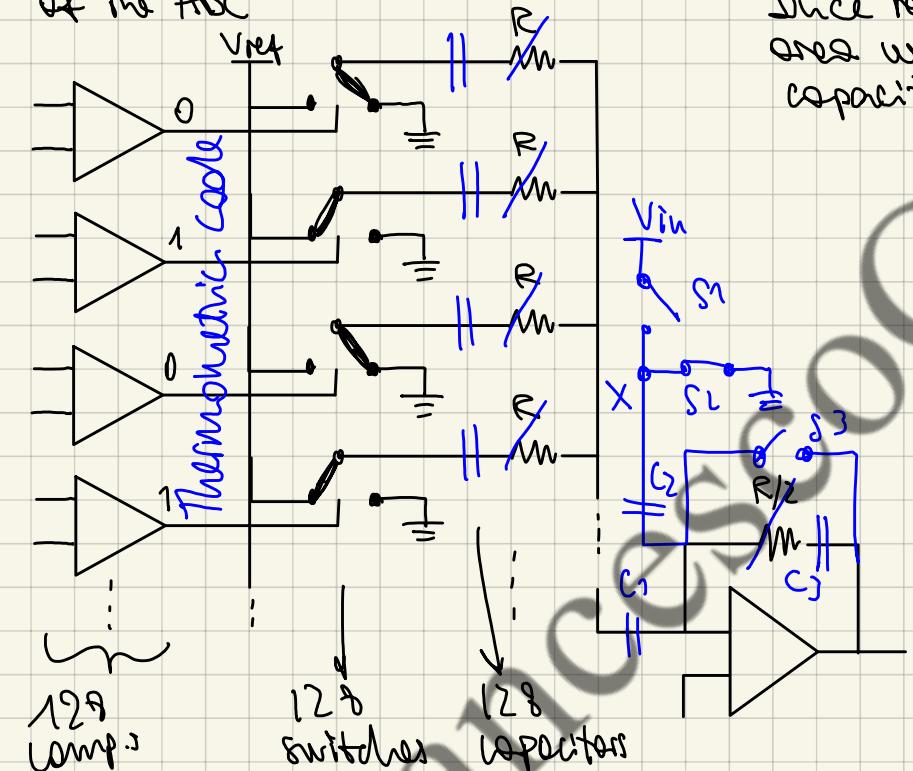
The comparators suffers of its own  $V_{os}$ , but we stored it across  $C_1$  so the two compensate each other

and we apply  $V_{il} = V_{ref}$  which will be compared to ground

so we get  $V_{il} - V_{ref}$  which will be compared to ground

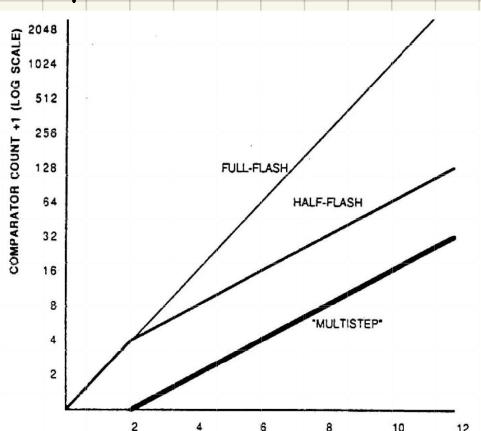
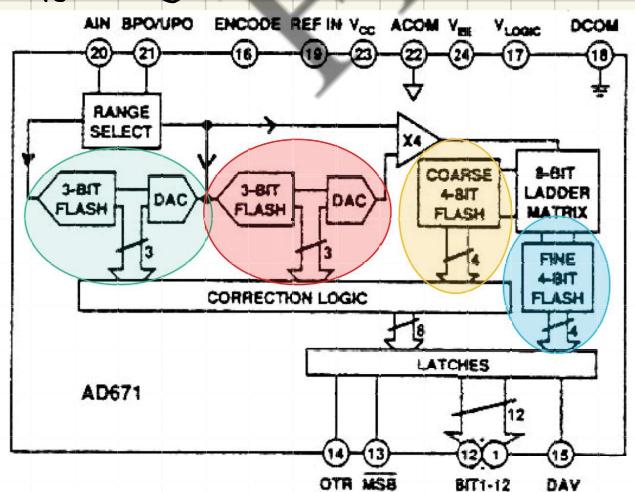


Comparators  
of the ADC



Since resistors require a huge amount of area we can substitute the resistors w/ capacitors

MULTISTEP FLASH ADC → we have a cascade of more than 2 flash ADCs

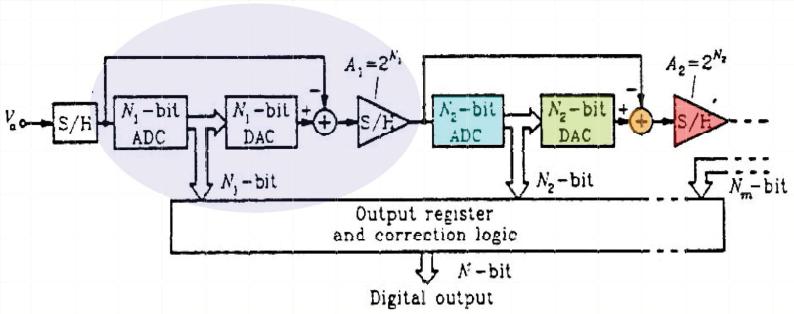


## Pipelined ADC

We've seen that the half-flash and the multi-step configurations allow, in respect of flash ADCs, to greatly minimize the area occupation, the power consumption and the conversion speed.

At the same time, many applications require high resolution and high sampling frequency.

⇒ The demand for high sampling freq, low area occupation and moderate power consumption is fully satisfied by using **PIPELINED CONVERTERS**.



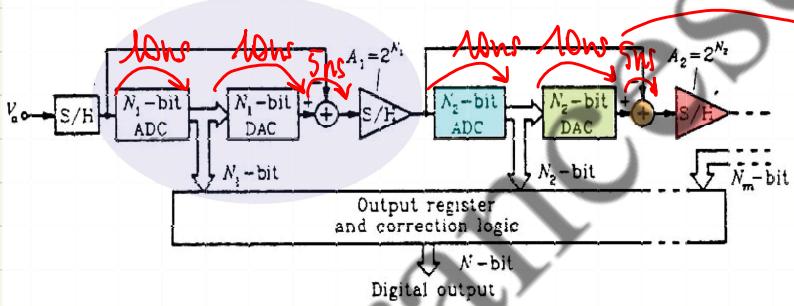
A typical **Pipelined ADC** is composed of  $m$  stages each of which contains:

- an S/H amplifier
- a low resolution flash ADC ( $1 \frac{1}{4}$  bit)
- a DAC
- an analog adder

From one stage to the other there always be a S/H, in this way after the conversion of the  $i$ -th stage, the residual which reaches the  $(i+1)$ -th stage is maintained by the S/H.

In doing so, the  $i$ -th stage can reach a new sample.

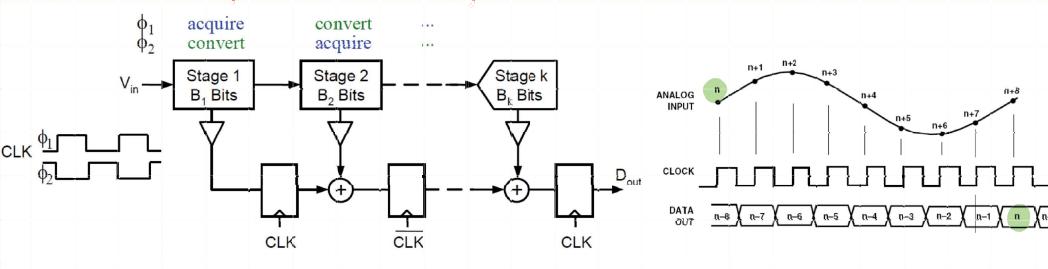
**NOTICE:** A pipelined architecture is a multistep ADC where each stage has sufficient time to convert, but it's not converting something that is propagating synchronously through the network, but each stage is converting something that is kept constant by the previous S/H.



in a 4 blocks ADC, one sample takes 100ns to convert  
(in a flash ADC it took just 10ns)  
but, at every clock pulse we can obtain my output but we can put together the result of the 1st coarse conversion w/ that one of the coarse conversion that happens 25ns later and so on

⇒ In a pipelined ADC w/  $m$  stages the sampling frequency is equal to the clock frequency, while every conversion requires  $m$  clock cycles, and therefore the code for a certain sample appears at the output after a certain delay, named **LATENCY** or **Pipeline Delay**.

## Pipelined ADC LATENCY



The output conversion of the  $n$ -th sample will be available after a given amount of clock cycles equal to the # of stages in which the ADC is divided

⇒ We need to properly choose the # of stages

ADVANTAGE: we've got an output every Tclock

ISSUE: LATENCY: Using 8 stages, the output related to the i-th sample will be available after 8 Tclock ( $= T_{conv}$ )

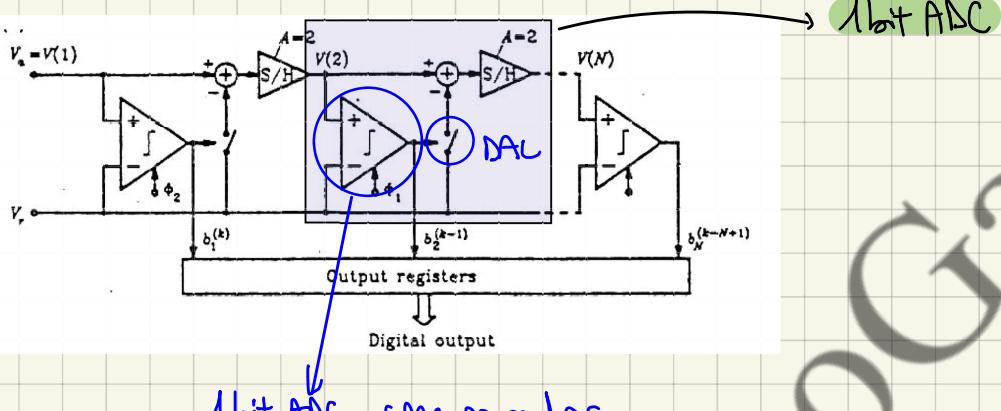
↓  
but it's not a real issue b/c we'll make the ADC operating in FREE-RUNNING mode

### "EXTREME" IMPLEMENTATION OF A PIPELINED ADC

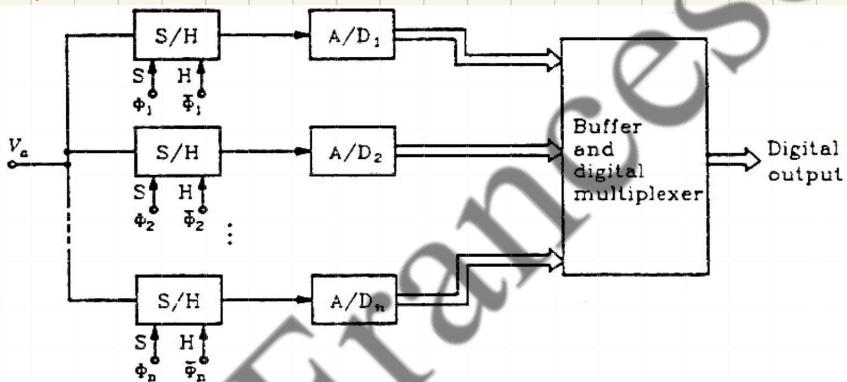
We wanna design a 12 bit ADC, instead of dividing them in 4+4+4, we use 12 stages of 1 bit each

What's a 1 bit ADC? → it's a comparator

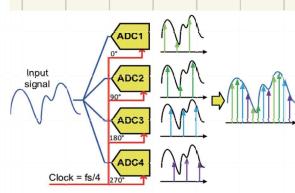
What's a 1 bit DAC? → it's a switch



### TIME-INTERLEAVED ADC



To obtain a high sampling freq. we can use two or more ADCs connected in parallel, obtaining a structure named the "TIME INTERLEAVED ARRAY"



This structure is made of C channels in each of which is present a S/H and a converter w/ N bits

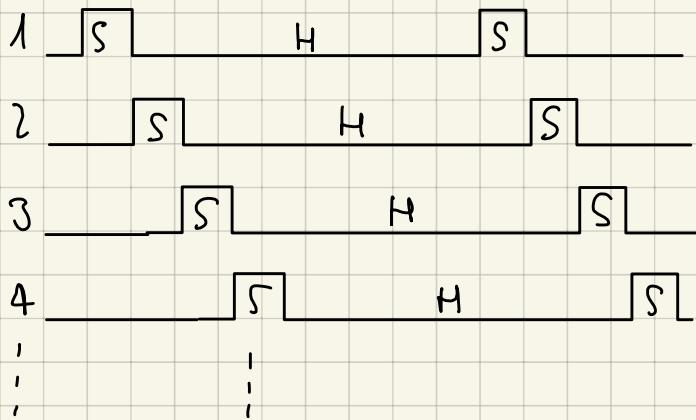
In every channel the input signal is sampled w/  $f_s = 1/C T_{clock}$

The channels work sequentially. Therefore the sampling frequency of the whole ADC is n times greater than one of the single converter whereas the resolution and the maximum frequency of the input signal are unchanged

example: if we want  $f_s = 16\text{Gps}$  w/ a 10-ADC it means that in principle  $T_{conv}$  must be  $T_{conv} = 1/f_s = 1\text{ns}$  which is impossible to obtain w/ a 10 bit ADC

→ we use 10 channels, each of them composed by a S/H and a 10-ADC w/  $T_{conv} = 10\text{ ns}$

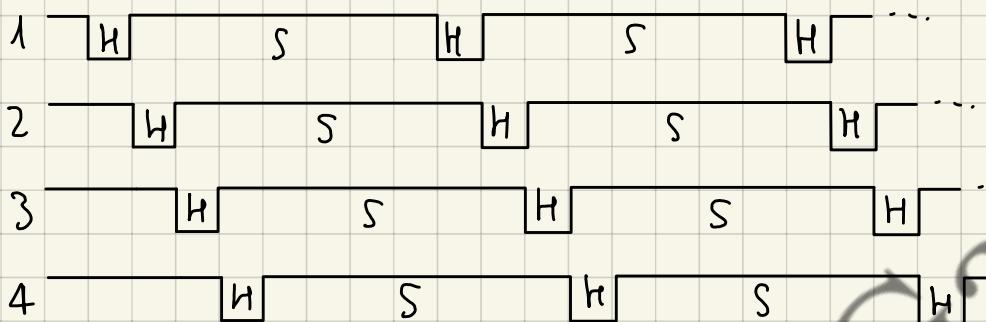
## ① POSSIBLE IMPLEMENTATION



→ This is a very stupid approach bcz during sample the switch is closed for a very short time and so the S/H should be very fast in acquiring the signal



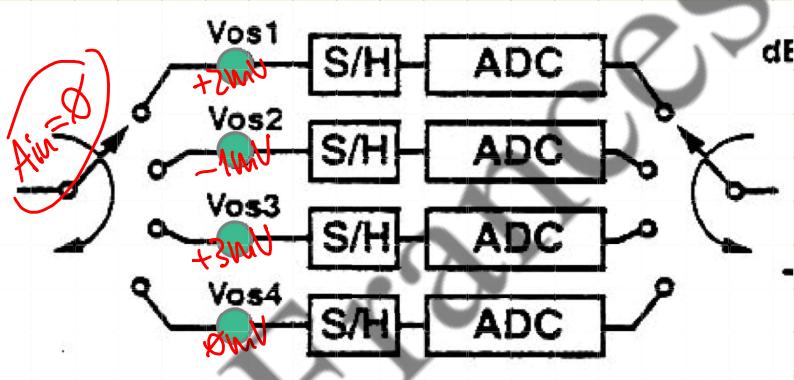
## ② POSSIBLE APPROACH → BETTER!!



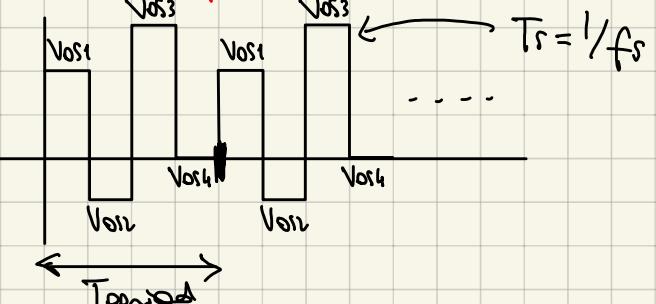
$$T_H \geq T_{conv, each ADC} = 10\text{ns}$$

### ERRORS:

#### ① OFFSETS

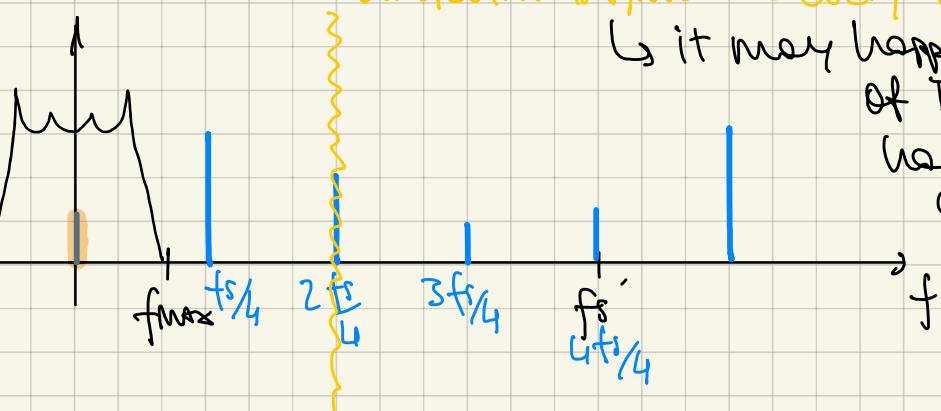


TYPICAL EXAMS'S QUESTION  
 What's the error that would you read out from the 4 ADCs?



$$\text{OUTPUT SIGNAL: } f = \frac{1}{T_{\text{period}}} = \frac{1}{4T_s} = \frac{1}{4} f_s$$

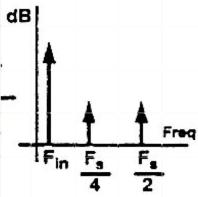
### OUTPUT SPECTRUM



our spectrum is symmetric every  $f_s/2$

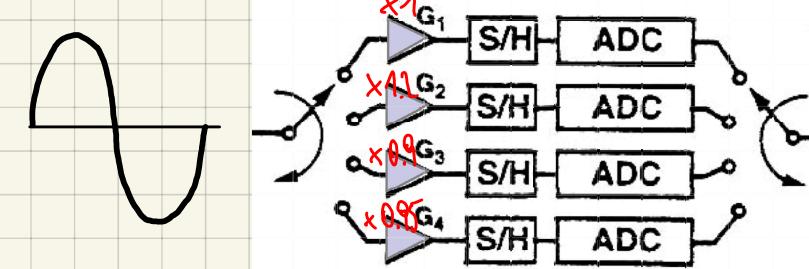
↳ it may happen that one of the following harmonic falls inside our baseband

⇒ An offset error in the ADC appears as a tone at  $f_s/C$ ,  $2f_s/C$ ,  $3f_s/C$ , ...



## ② GAINS

Imagine that every S/H has its own gain



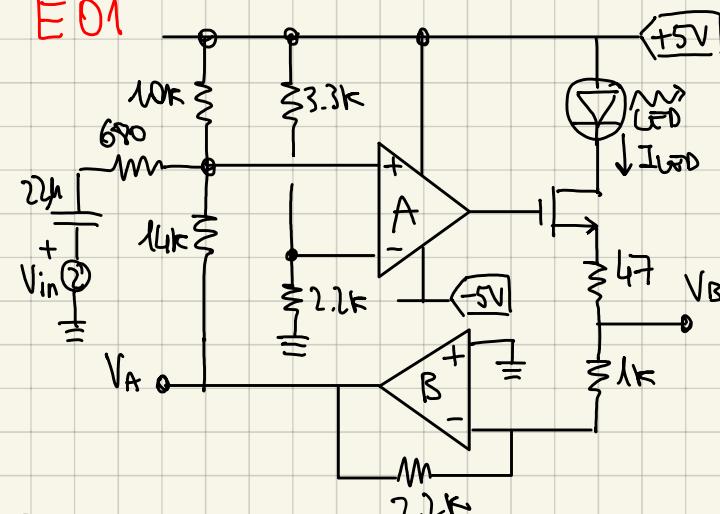
TO FINISH!!!

Francesco Gavetti

# ES15 - EXERCISES

01/12/2021

E01



OpAmps:  $A_v = 50 \text{V/mV}$

$$V_{os} = 5 \text{mV}$$

$$I_B = 200 \mu\text{A}$$

$$\text{MOSFET: } k = \frac{1}{2} \mu_n C_{ox} = 10 \text{mA/V}^2$$

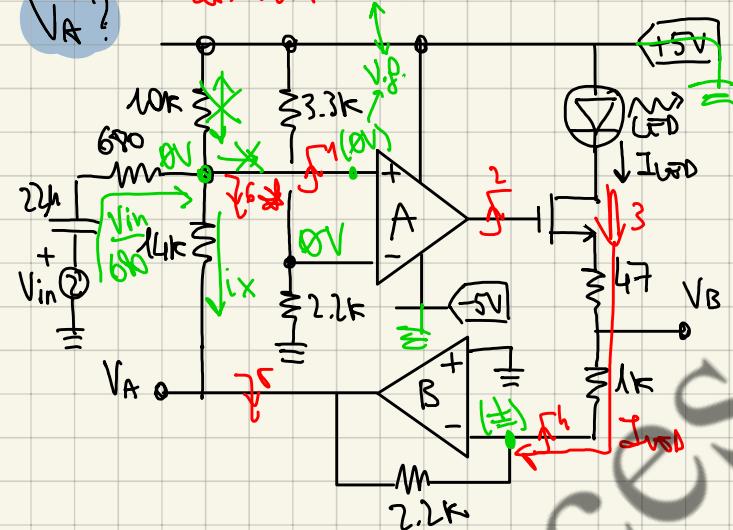
$$V_T = 0.8 \text{V}$$

(a) Compute the relationships of  $V_A$ ,  $V_B$  and  $I_{LED}$  vs.  $V_{in}$ .

## SMALL SIGNAL ANALYSIS

\*NEGATIVE FB

$V_A$ ?



The Opamp A has an amplification of  $A_v$  while the Opamp B has an amplification of  $(-2.2k)$  but it's an inverting config.

we proceed in  
anti-clockwise direction  
signal

When we study  $V_{in}$  we have to turn off all the other sources, included the PS. To perform the small signal analysis

$$@DC \rightarrow C_{open} \rightarrow \frac{V_A(0)}{V_{in}} = \infty$$

$$@HF \rightarrow C_{short} \rightarrow V_A = -V_{in} \frac{14k}{680} \Rightarrow \frac{V_A(\infty)}{V_{in}} = -20.6 \frac{\text{V}}{\text{V}}$$

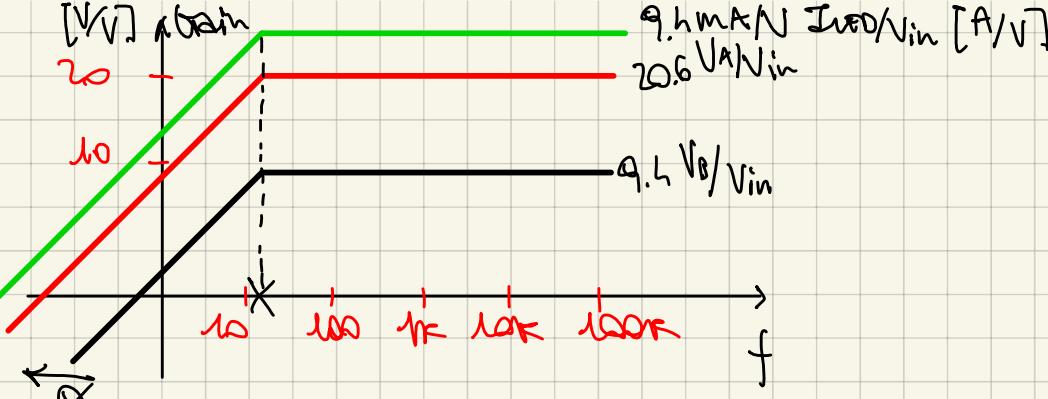
$V_B$ ?

$$@DC \rightarrow \frac{V_B(0)}{V_{in}} = \infty$$

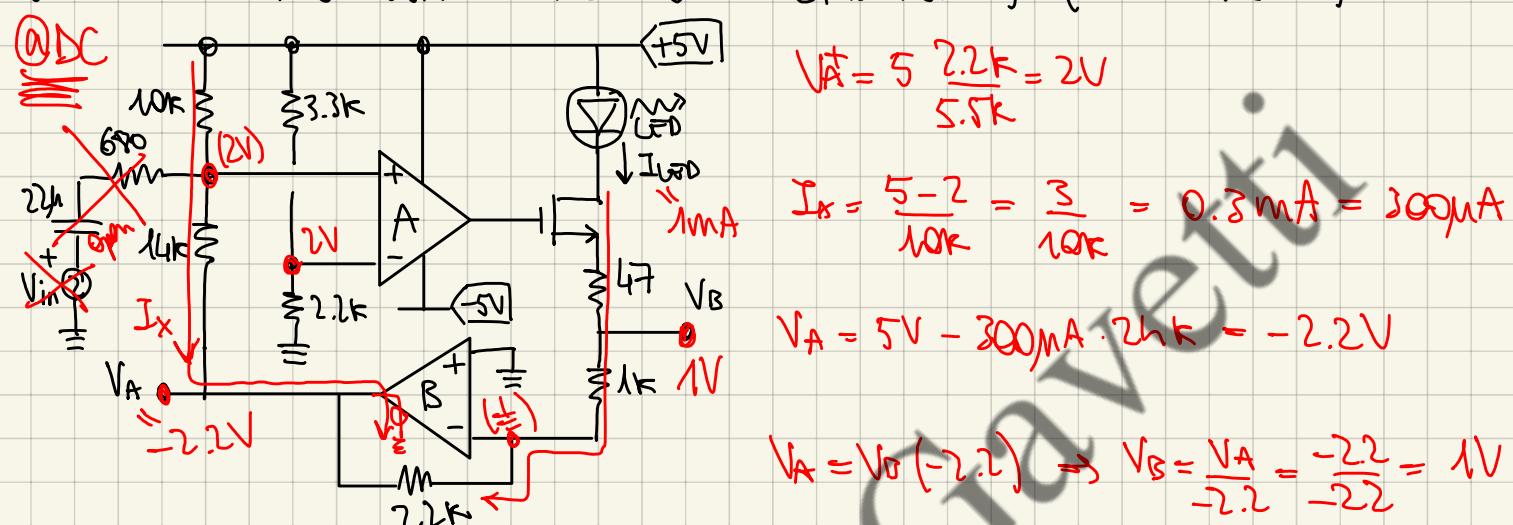
$$@HF \rightarrow V_A = V_B \left( -\frac{2.2k}{1k} \right) \rightarrow -20.6 V_{in} = V_B (-2.2) \Rightarrow \frac{V_B(\infty)}{V_{in}} \approx 9.4 \frac{\text{V}}{\text{V}}$$

$I_{LED}$ ?

$$I_{LED} = \frac{V_B}{1k} \Rightarrow \frac{I_{LED}(\infty)}{V_{in}} = \frac{9.4}{1k} = 9.4 \text{ mA/V}$$



LET'S STUDY THE POLARIZATION (THE BIAS POINT) (when  $V_{in} = \infty$ )



$$I_{LED} = \frac{V_S}{1k} = \frac{1V}{1k\Omega} = 1mA$$

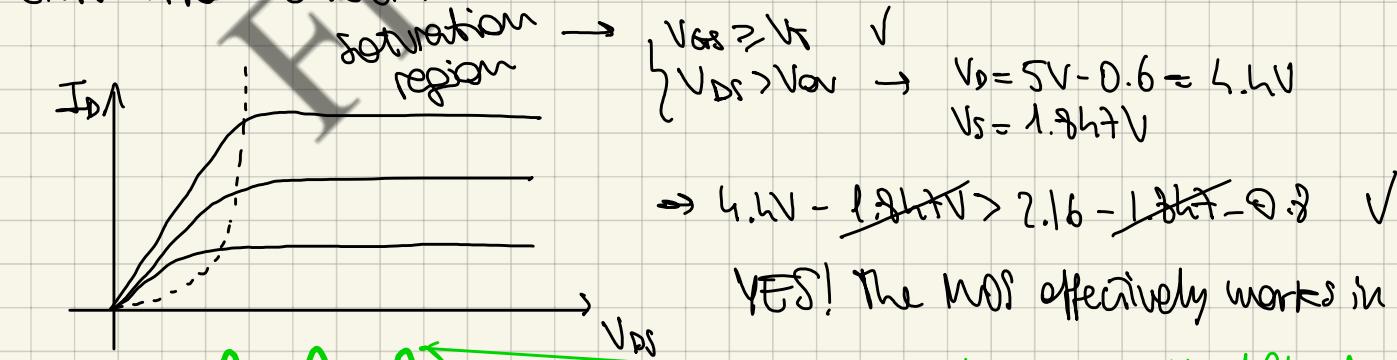
$$V_S = V_B + I_{LED} \cdot R_D = 1V + 47mA = 1.047V \quad "0.316"$$

$$i_d = k (V_{DS} - V_T)^2 \rightarrow V_G = V_S + \sqrt{\frac{i_d}{k}} = 1.047V + \sqrt{\frac{1mA}{100mA/V^2}} = 2.16V$$

$$\left\{ \begin{array}{l} V_A = -2.2V \\ V_B = 1V \\ I_{LED} = 1mA \end{array} \right.$$

POLARIZATION  
VAVES

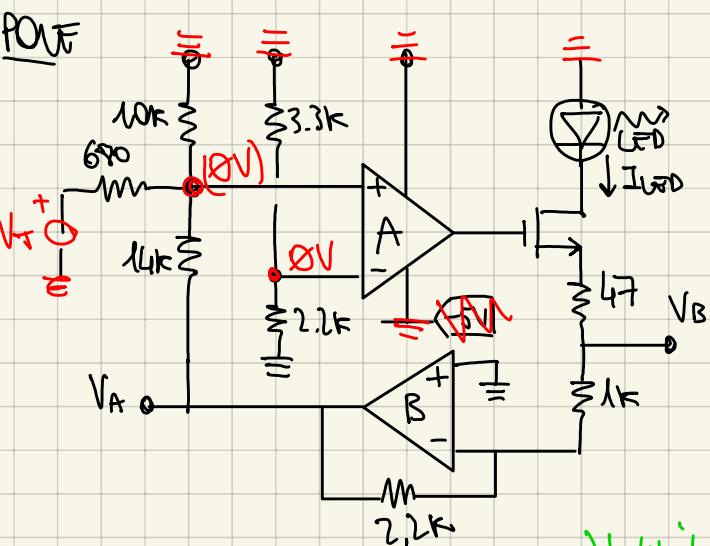
SATURATION CHECK:



$$I_{D,\max} = K \left( \sqrt{V_{GS,\max} - V_T} \right)^2 = 10 \frac{\text{mA}}{\sqrt{\text{V}}} |5\text{V} - 0.8\text{V}|^2 = 176 \text{mA}$$

$$I_{D,\min} = 0 \text{A} \quad (\text{bcz cannot invert})$$

(B) Compute The input pole and plot  $V_A(t)$  and  $I_{LED}$  waveforms when the input is a high frequency  $200\text{mV}_R$  sinusoid



$$f_p = \frac{1}{2\pi R C} = \frac{1}{2\pi \cdot 22\mu\text{F} \cdot 600\Omega} = 10.6 \text{ Hz}$$

Beware:

Small signal

$$100\text{mV} \cdot \frac{V_A(t)}{V_{in}} (0) = 2.06 \text{V}$$

DC contrib  
(no variation in time)

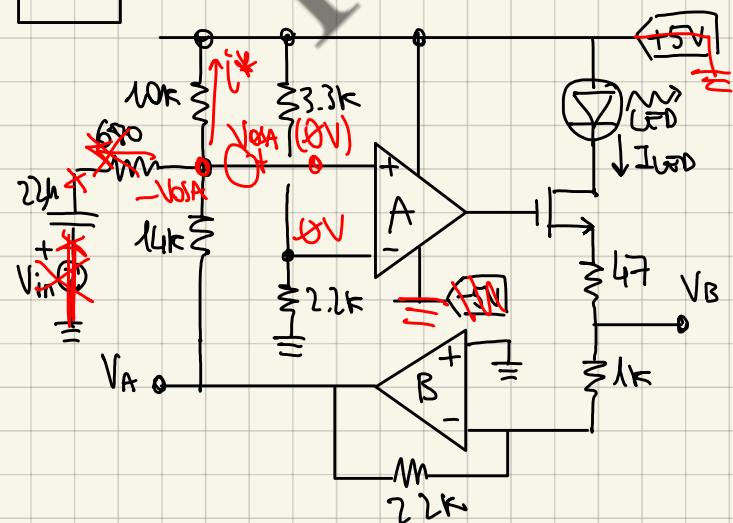
$V_{in}(t)$

$V_A(t)$  it stays to zero bcz  
 $V_A(\infty) = 0$



(C) Compute The maximum  $V_A$  static error due to  $V_{os}$  and  $I_{os}$  of The Opamp

$V_{osA}$



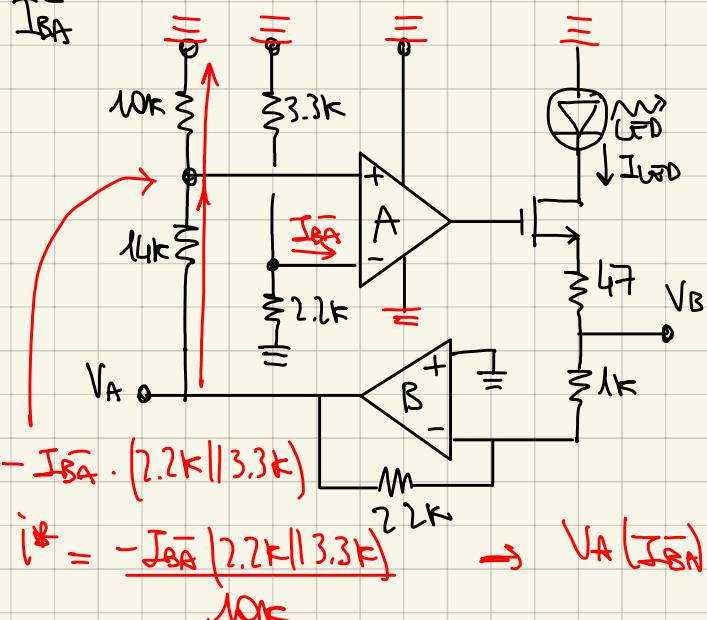
$$i^* = \frac{V_{osA}}{10\text{k}}$$

$$V_{Aos} = \frac{V_{osA}}{10\text{k}} \cdot 2.2\text{k} = \pm 5\text{mA} \cdot 2.4 =$$

$$= \pm 12\text{mV}$$

$$V_{osos} = \frac{\pm 12\text{mV}}{-2.2} = \mp 5.5\text{mV}$$

$$I_{DOS,os} = \frac{\mp 5.5\text{mV}}{1\text{k}} = \mp 5.5\text{nA}$$

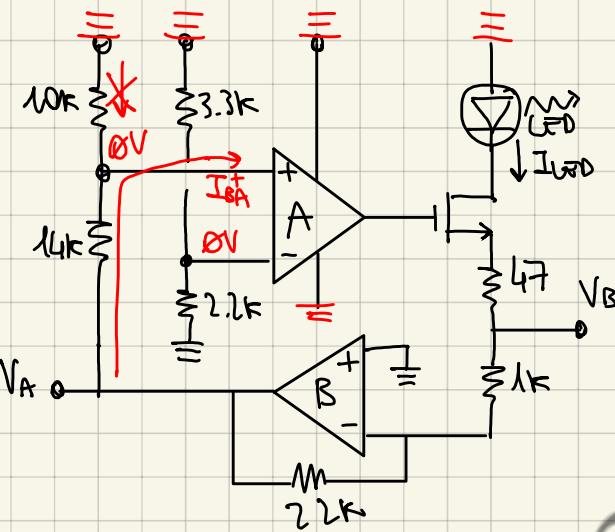
$I_{BA}^-$ 

$$-I_{BA}^- \cdot (2.2k \parallel 3.3k)$$

$$i^* = -\frac{I_{BA}^- (2.2k \parallel 3.3k)}{10k}$$

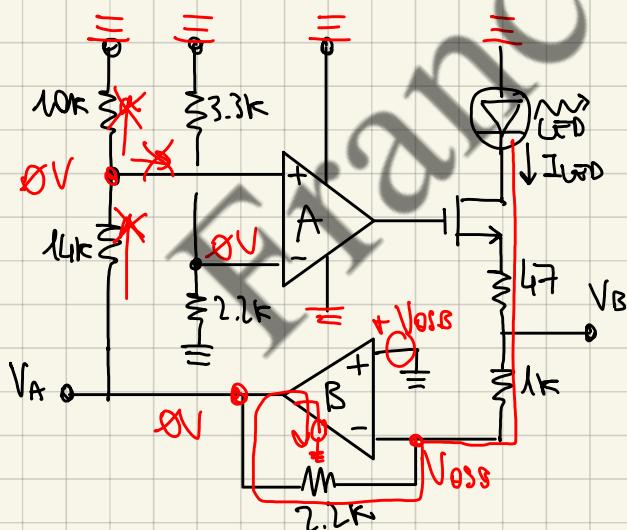
$$V_A (I_{BA}^-) = -I_{BA}^- (2.2k \parallel 3.3k) \cdot 2.2k =$$

$$= -0.63 \mu V \text{ negligible wrt } \pm 12mV$$

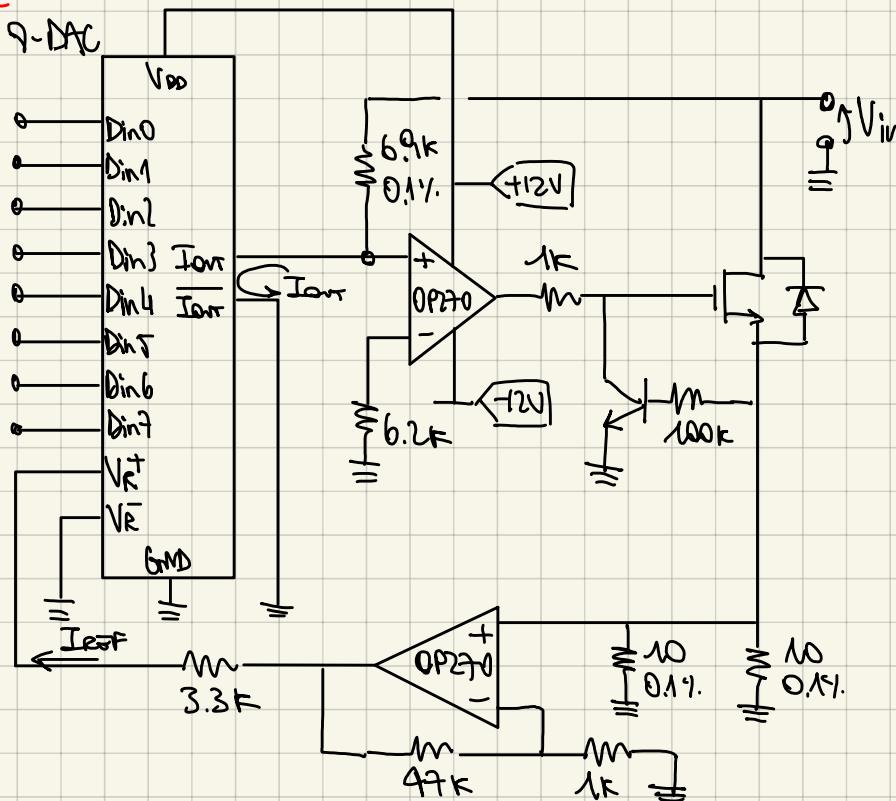
 $I_{BA}^+$ 

$$V_A (I_{BA}^+) = I_{BA}^+ \cdot 1k\Omega = 200 \mu A \cdot 1k\Omega = 2.8 \mu V$$

negligible



E02



## DAC providers:

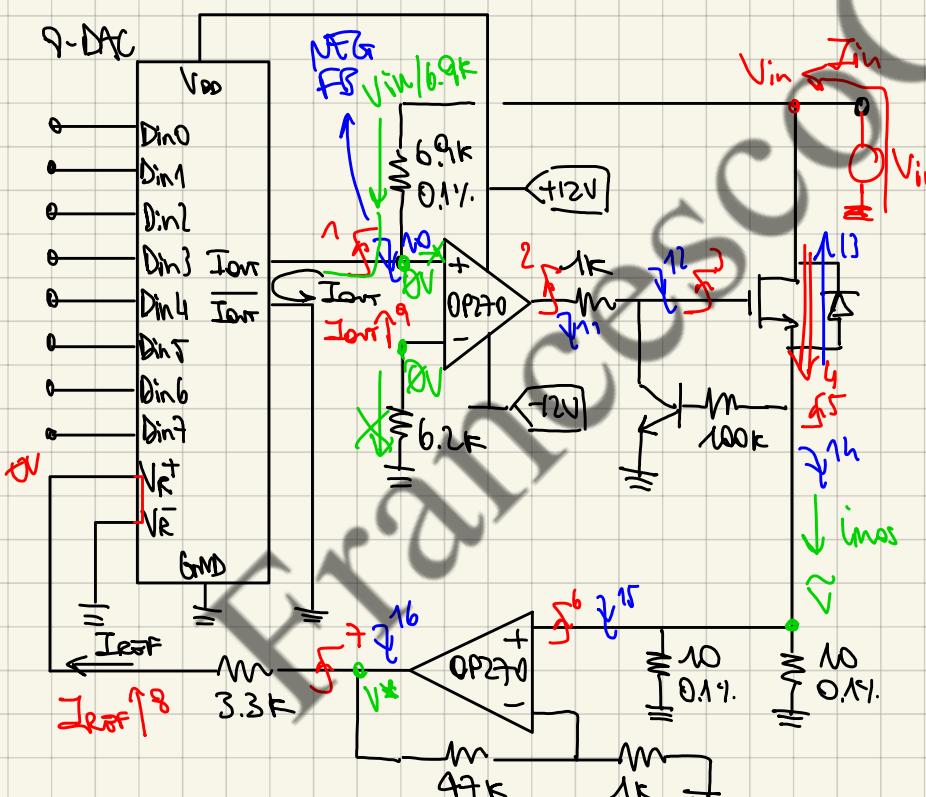
$$I_{out} = I_{ref} \cdot D_{in} / 256$$

and  $v_f$  @  $V_E^t$  input

$$\text{Opamps: } A_0 = 1500 \text{V/mV} = 1.5 \text{M}$$

$$GFWP = 5M\mu_t$$

Ⓐ Obtain The analytical relationship  $V_{in}/I_{in}$ , Or  $Q$  function of  $N$



$$I_{out} = \frac{V_{in}}{6.9k} = I_{ref} \frac{D_{in}}{?}$$

$$I_{ref} = \frac{V_{in}}{6.9k} \cdot \frac{2^8}{D_{in}} = \frac{V_{in}}{D_{in}} \cdot \frac{1}{2^8 k}$$

$$V^* = \text{Inf. } 3.3\pi$$

$$V^* = \sum \left( 1 + \frac{t}{k} \right)^k$$

$$Z = \frac{V^*}{47} = \frac{V_{in}}{Dm} \frac{33k}{27 \cdot 47} = \frac{V_{in}}{Dm} 2.5$$

$$I_{\text{mer}} = \frac{2}{10 \text{Hg}} = \frac{2}{5} = \frac{\text{Vin}}{\text{R}_\text{in}} \cdot \frac{1}{2}$$

$$\Rightarrow \frac{\sqrt{I_2}}{I_2} = 2\zeta \cdot D_{in} = R_{eq}$$

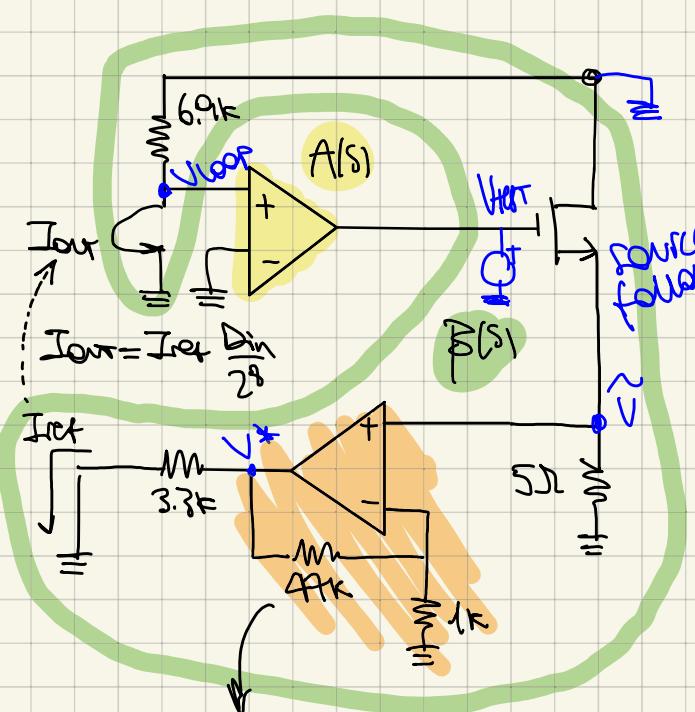
$$D_m = 0 \quad \rightarrow \quad R_{\text{eff}} = 83$$

$$P_{\text{obs}} = ?$$

$$D_{\text{h}} = 12.8 \rightarrow R_{\text{eq}} = 7.56 \text{ cm}$$

$$D_{in} = 255 \rightarrow Reg = 5101$$

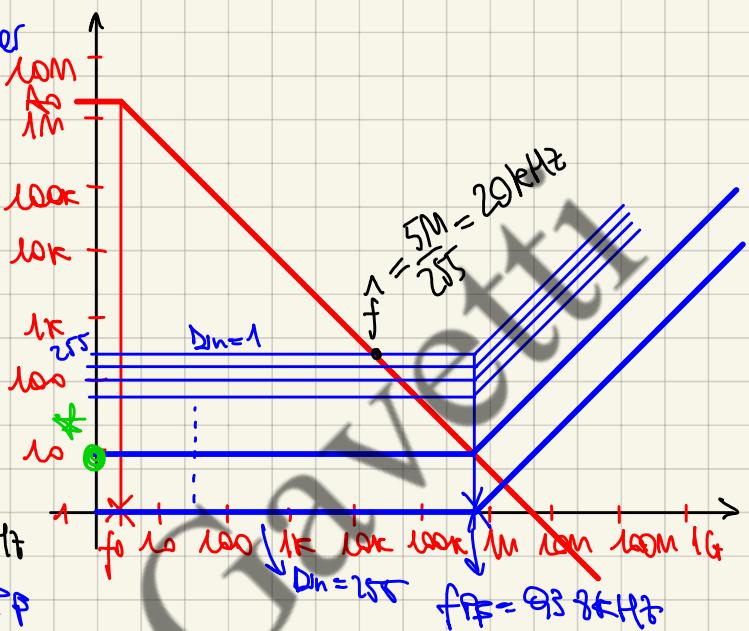
(B) Reckon if the stage is stable or not when  $Din = 255$ . Moreover, tell if mobility improves by reducing  $Din$  (Hint: assume  $\gamma_{gm,NOI} = 1.25 \text{ nA}$  and ignore the role of the BJT)



$$G_{FBWP} = 5M, G_1 = 48, f^* = \frac{5M}{48} = 104 \text{ kHz}$$

$$A_o = 1.5M = 20 \log_{10}(1.5 \cdot 10^6) = 123 \text{ dB}$$

$$f_o = \frac{G_{FBWP}}{A_o} = \frac{5M}{1.5M} = 3.3 \text{ Hz}$$



$$\beta = \frac{5M}{5M + \frac{1}{g_m}} \cdot \left(1 + \frac{48k}{1k}\right) \cdot \frac{1}{3.3k} \cdot \frac{Din}{2^8} \cdot (-6.9k) = -\frac{Din}{255}$$

$$* \frac{1}{\beta_x} \cdot 938 \text{ kHz} = 5 \text{ MHz}$$

$$\frac{1}{\beta_x} = \frac{5M}{938k} = 48 \leq Din \text{ stable}$$

if  $Din < 48 \Rightarrow \text{UNSTABLE}$

if  $Din = 48$  (20-40 degree angle)

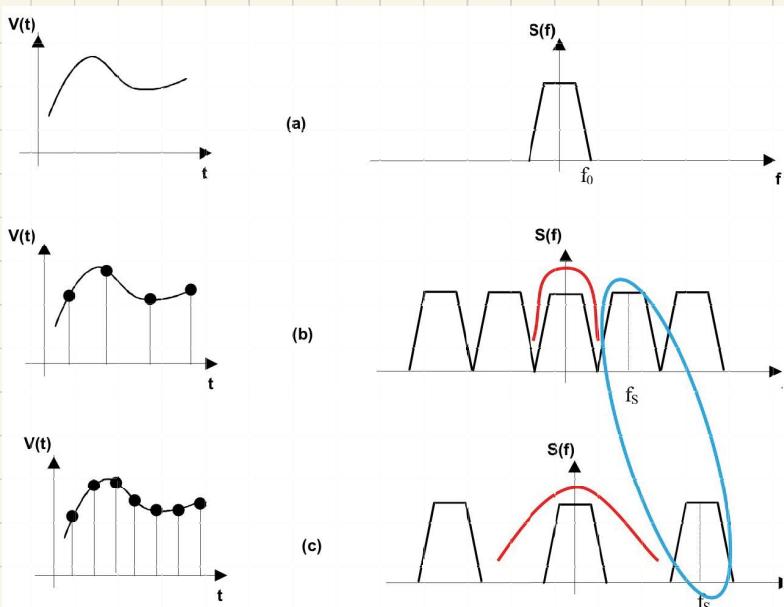
$$\text{PM} = 45^\circ$$

# ES15 - OVERSAMPLING AND SIGMA-DELTA

14/12/2021

- why increasing  $f_s$  beyond the Shannon minimum
- standard oversampling
- advanced sigma-delta oversampling
- why 1bit is equal or even better than many?
- need for DSP for digital filtering and decimation

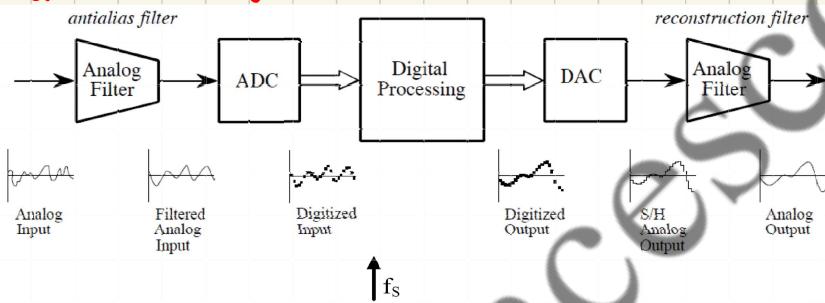
## SAMPLING



if  $T_s$  is long  $\rightarrow f_s$  is low and so the replicas will be very close to each other or they can even overlap if  $f_s < 2f_{\text{signal}}$   
Furthermore if we have to pass only the baseband and kill all the other replicas our LPF must be very selective

↳ decreasing  $T_s$ , we increase the sampling frequency  $f_s$ , so the replicas will be more apart from each other and our LPF can be more relaxed (lower # of poles)

## QUANTIZATION



The analog filter at the input and the one at the output must be identical but they should let the signal pass up to  $f_s$  killing all the other replicas

Due to the fact that the ADC provides a quantization of the vertical amplitude w/ steps equal to 1LSB, the ADC introduces an error into the signal, but the signal that feed the DSP is not equal to the original one but it will be affected by the so called QUANTIZATION ERROR

$$\overline{v_q^2} = \left( \frac{FSR}{2^n \sqrt{12}} \right)^2 = \frac{LSB^2}{12}$$

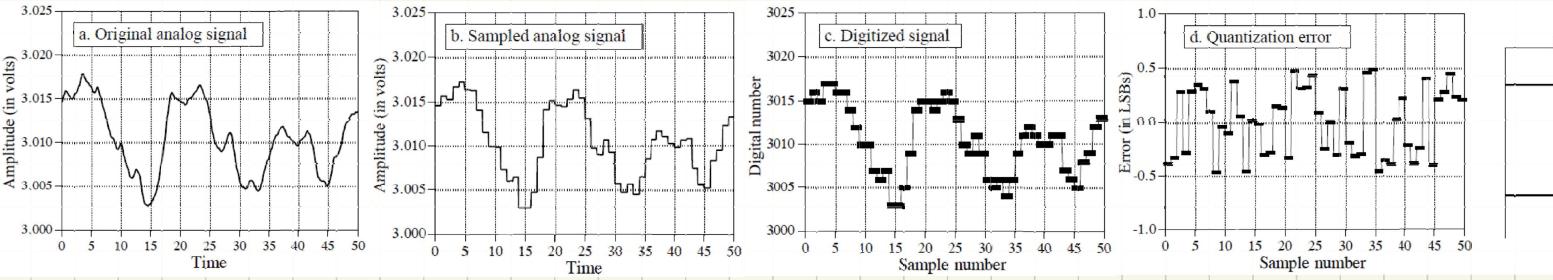
POWER OF QUANTIZATION ERROR

$$Q_q(f) = \frac{\overline{v_q^2}}{f_s/2}$$

QUANTIZATION ERROR'S SPECTRAL DENSITY

↳ The spectrum is symmetric respect to  $f_s/2$

example: FSR=5V  $f_s=10\text{kHz}$   $\overline{v_q^2}=(352\mu\text{V})^2$   $(Q_q(f)=(3.52\mu\text{V}/\sqrt{\text{Hz}}))^2$

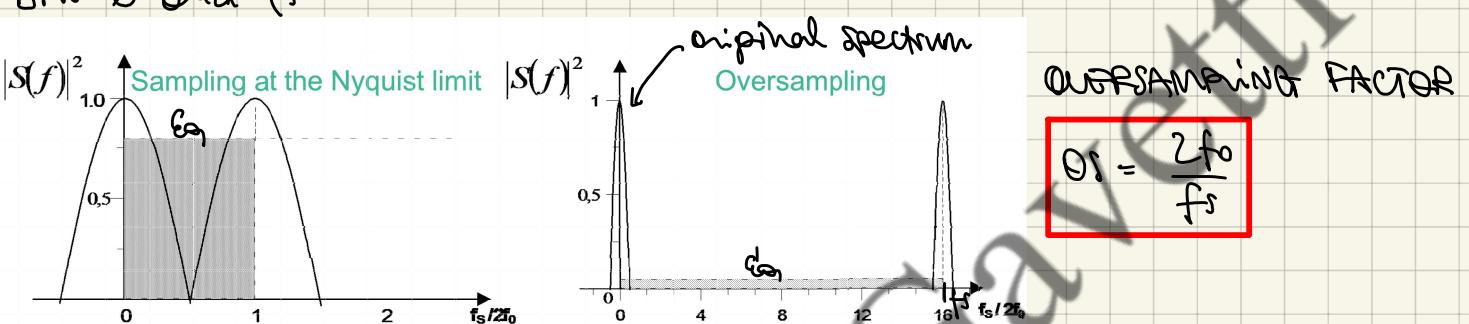


$$\text{Quantization error} = \text{Sampled signal} - \text{Digitized signal} : d = b - c$$

If the ADC is well designed  $|E_Q| \leq \frac{1}{2} \text{ LSB}$

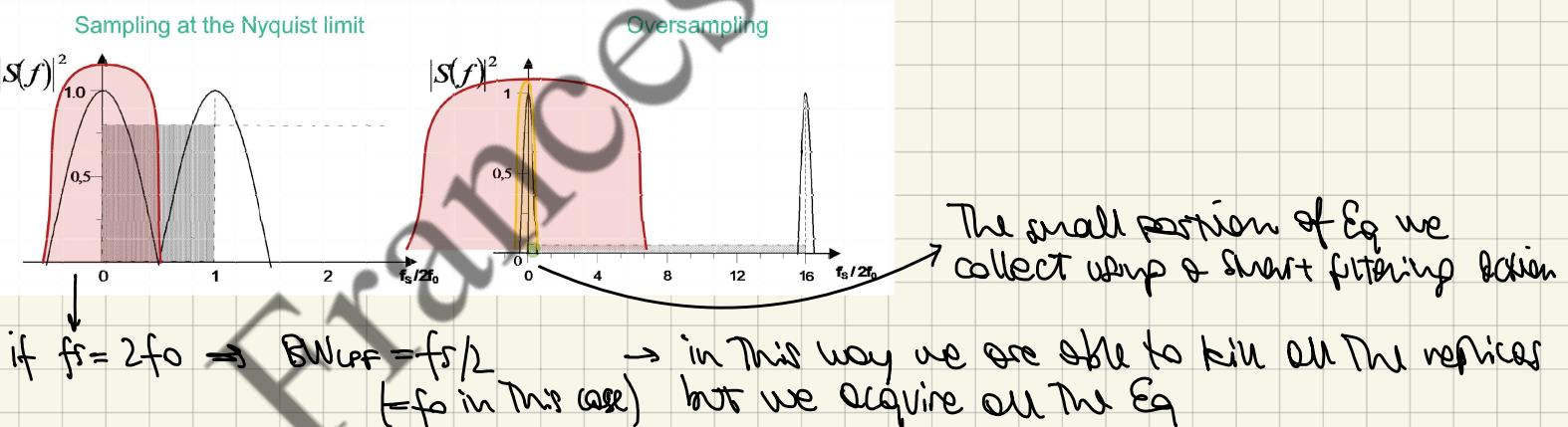
## OVERSAMPLING

If we perform our sampling @ the Shannon/Nyquist limit  $f_S = 2f_{\max}$  we end up that the replicas touch @  $f_{\max}$  and  $E_Q$  is distributed uniformly b/w 0 and  $f_S$



Otherwise, if we employ  $f_S > 2f_{\max}$  (e.g.  $f_S = 16(2f_{\max})$ ) The two replicas are very separated one another and  $E_Q$  gets spread over a wider frequency range

$$E'_Q = \frac{E_Q}{16} \rightarrow \text{we can get an improvement if we perform a smart filtering action}$$



**OVERSAMPLING**: even in this case we can use a LPF w/  $BW = f_S/2$  but again we reconstruct our original spectrum, killing all the replicas, but we also recollect all the  $E_Q$  the ADC had introduced

**SMART FILTERING**: let's filtering up to  $f_0$  and not  $f_S/2$ .

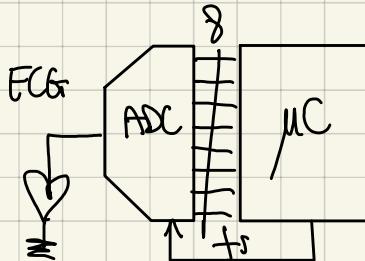
In this way we are able to fully reconstruct our original signal, but we'll recollect only a small portion of  $E_Q$

↓  
we improve The SNR !!!

CLASSICAL ADC :  $\text{SNR}_{\text{th}} = \frac{V_{\text{eff}}}{q} = \frac{\text{FSR}/2\sqrt{2}}{\Delta/\sqrt{2}} = \frac{\text{FSR}/2\sqrt{2}}{\text{FSR}/2\sqrt{n}\sqrt{2}} = 2^{n-1} \sqrt{6} \approx 6.02 n + 1.76$

OVERSAMPED ADC :  $\text{SNR} = \frac{V_{\text{eff}}^2}{\text{FSR}^2/16S} = \text{SNR}_{\text{th}} \cdot \sqrt{OS}$

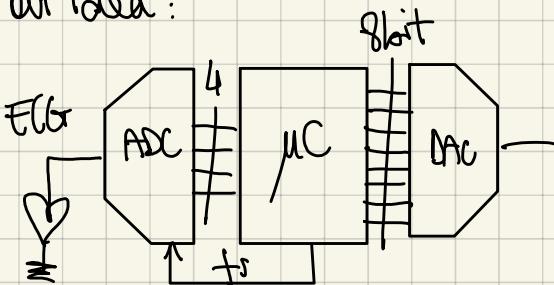
RESULT : we get an improvement of  $\sqrt{2} = 3 \text{dB} = \frac{1}{2} \text{bit}$  every doubling of  $f_s$ !



$$\text{SNR}_{\text{ideal}} = 6.02 \cdot 8 + 1.76 \approx 50 \text{dB}$$

$$f_{\text{max}} \leq 250 \text{Hz} \quad f_s = 2f_{\text{max}} = 500 \text{sps} \quad T_s = 1/f_s = 2 \text{ms}$$

Our idea :

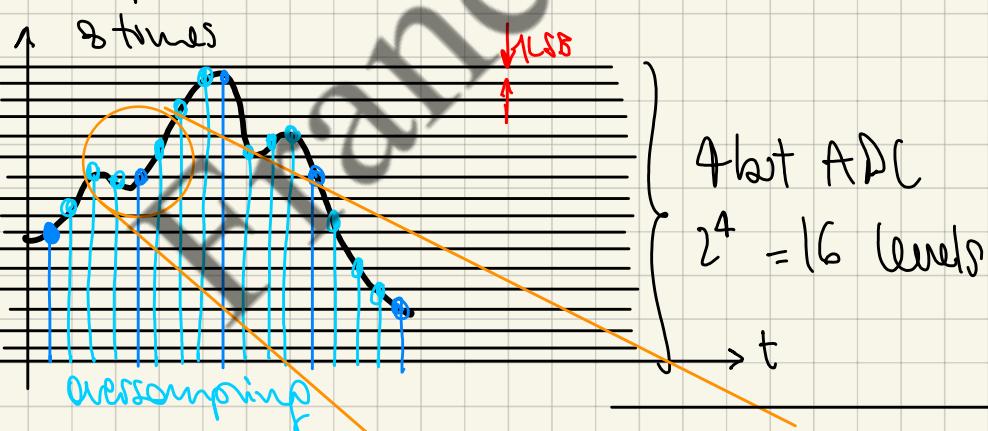


We want our μC to be able to provide 8 bit in terms of signal quality at the output

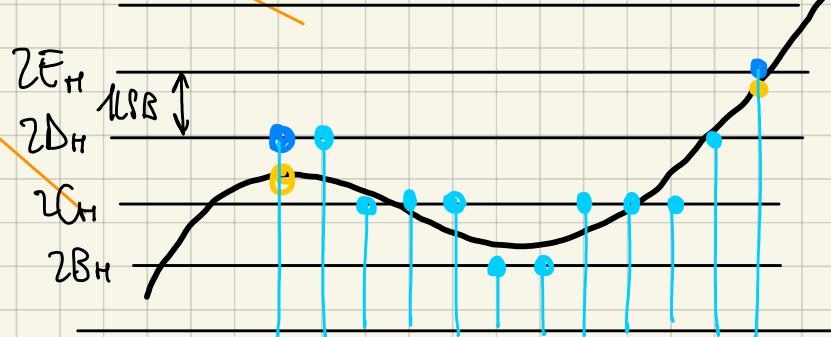
$$f_{\text{max}} \leq 250 \text{Hz} \quad \text{Unit -4 bit} = 4 \text{bit} \rightarrow \text{improvement needed}$$

Since we get a  $\frac{1}{2}$  bit improvement every doubling of  $f_s$ , if we want to achieve a 4 bit improvement we need to implement 8 doublings:

$$OS = 2 \cdot 2 \cdot 2 \cdots 2 = 2^8 = 256 \Rightarrow f_s = 2f_{\text{max}} \cdot OS = 256 (500 \text{sps}) = 128 \text{kspss}$$

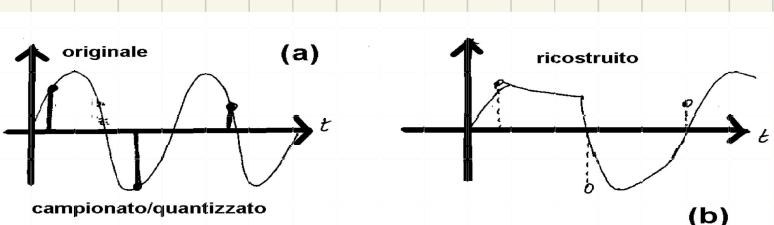


Oversampling doesn't give us more levels



How the μC can achieve an improvement in terms of its own resolution?

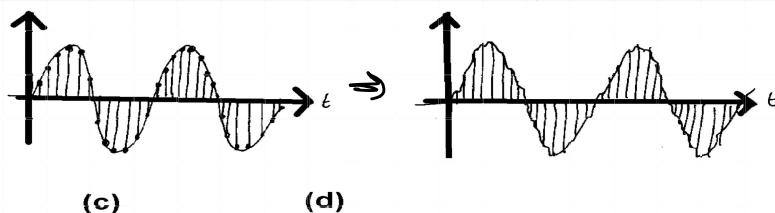
The improvement can be achieved by filtering



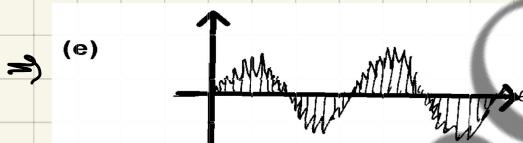
Using the Nyquist frequency and a LPF the quality of the reconstructed signal is poor

If we oversample the signal we have 2 possibilities:

- we use the same filter we'd use @ the Nyquist sampling rate

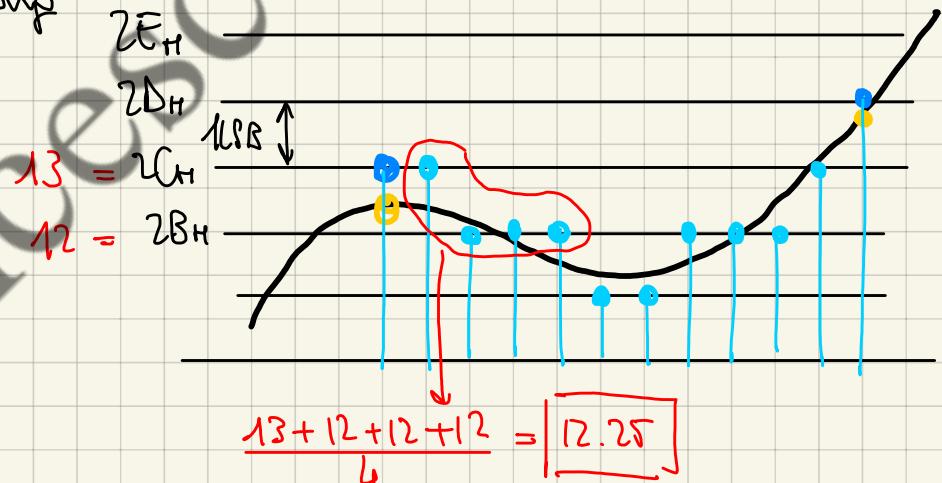
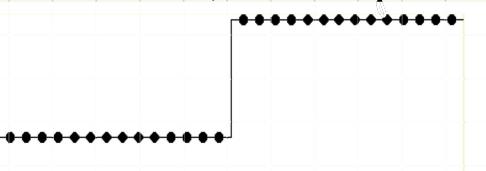


- we use a filter @  $f_s/2$ , which means that it has an higher bandpass so the T is much shorter, this means that the filter is much faster in following all the new samples



## DIGITAL FILTERING

The idea is performing averaging



After averaging our μC will have the possibility of having an arithmetics which is no more limited to the 7 bit

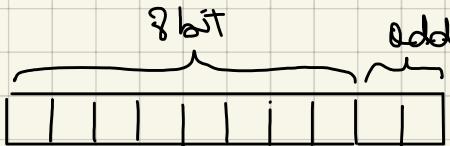


Averaging allows us to have a higher resolution in the μC's arithmetics



This means that the accumulator of the μC should be able to perform sum btw different samples and then the division by the # of samples considered

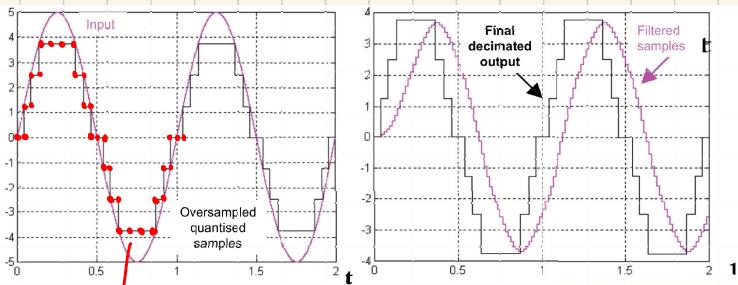
NOTICE: The division by 4 means that we have to perform a shift leftwards of the accumulator register by 2 bit  
→ a shift on the left of 1 bit corresponds to a division by 2



Added resolution due to the decimal part resulting from  
the filtering action

⇒ oversampling provides an improvement in the # of levels inside the μL  
Thanks To The digital oversampling of samples

EXAMPLE:  $f_0 = 1\text{kHz}$   $\theta S = 16$   $f_S = 32\text{fps}$   $n = 3\text{bit}$

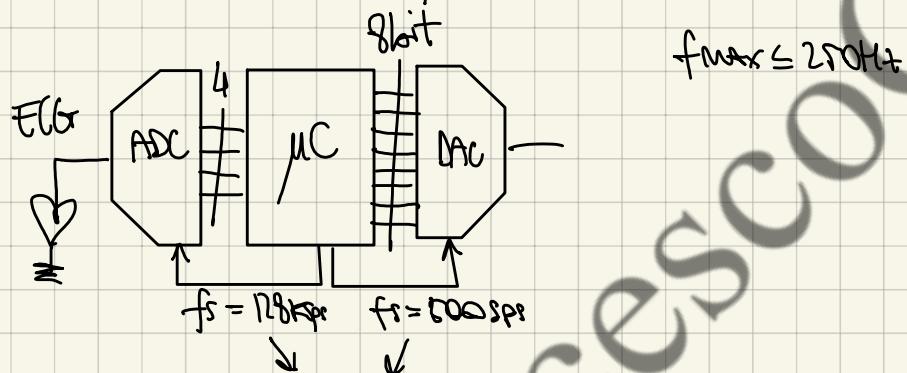


We acquire 16 samples more than what predicted by Shannon

even if we perform DS we're still limited by the resolution of the ADC (3-bit ADC)

Then filtering out the signal we basically perform an averaging which improves the resolution

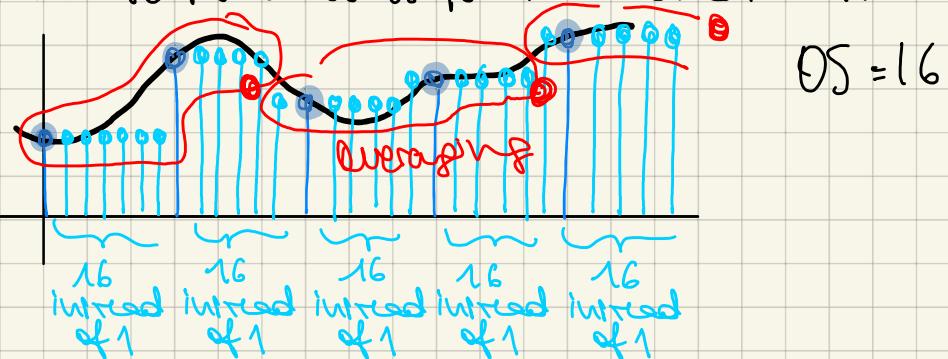
Sooner or later the PIC will provide the compressed stream to a DAC



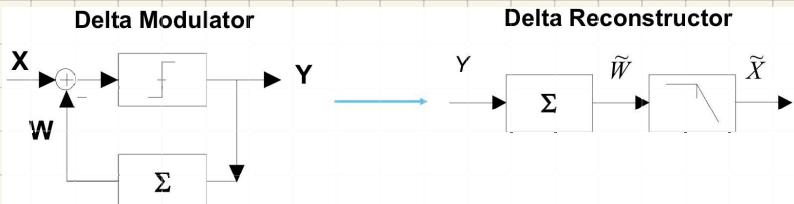
The two sampling freqs differ one each other: The  $f_s$  @ which the ADC makes the ADC operating is much higher than the Shannon limit in order to perform oversampling, while the DAC @ the output unit operate @  $f_s = 2f_{MAX} = 500\text{ s/s}$

At the output DAC we can perform the so called **DECIMATION** which means that once we acquired a larger # of samples, but we're operating the ADC in a oversampling way, Once we performed digital filtering, but we operated oversampling, Then we can simply store in our FLC only one sample out of n

- Select OS to provide the required bit improvement
  - ↓
  - The use the same OS for the decimation! !



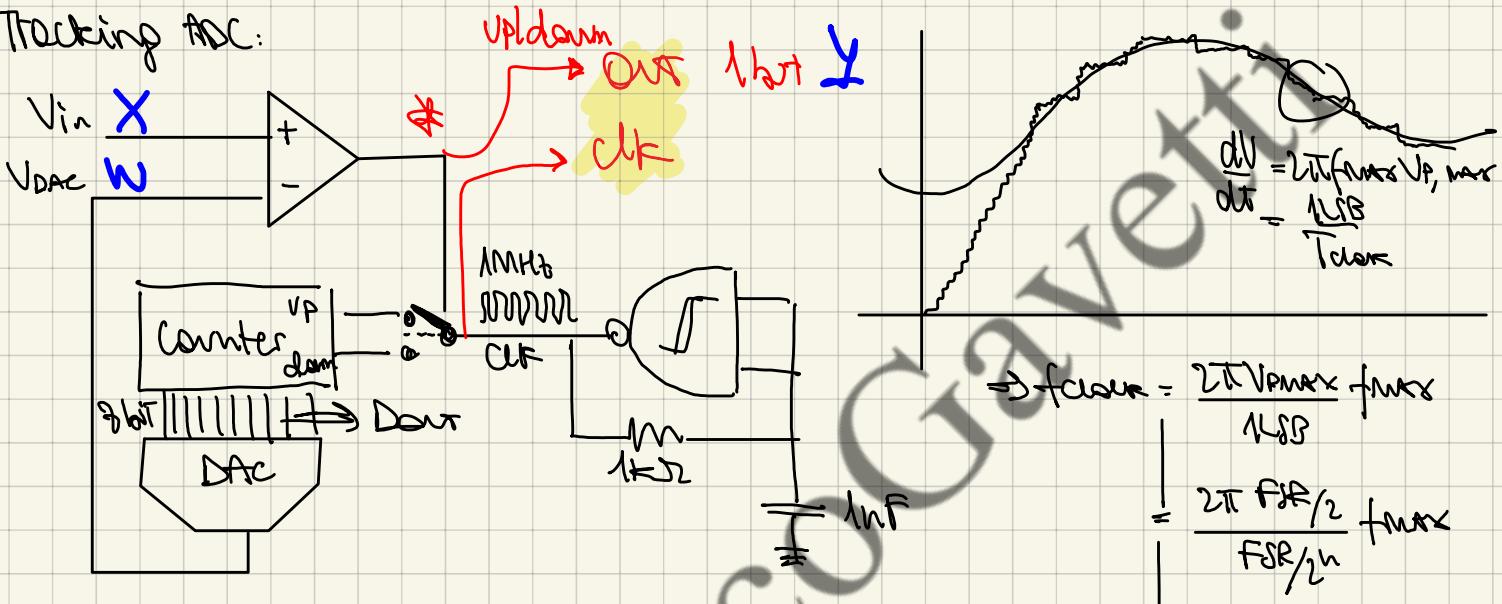
## DELTAMODULATION



We already discussed the STOICORE ADC which is able to track the original signal and then we improved it by allowing the counter to go both up and down, discovering the so called TRACKING ADC

The tracking ADC can also be described as a DELTA MODULATOR

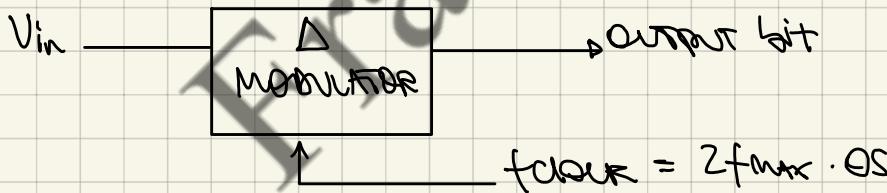
Tracking ADC:



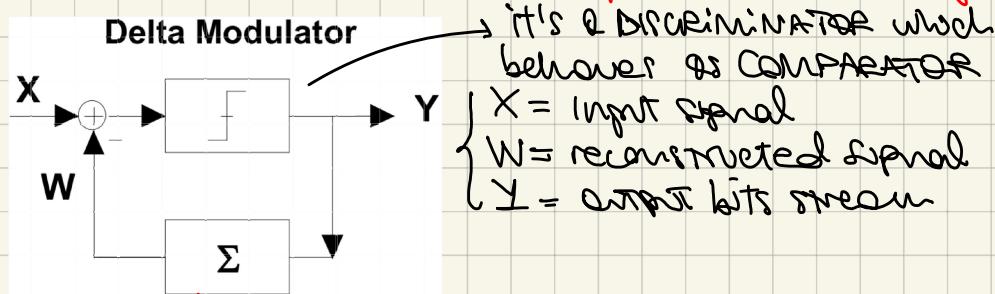
We get 1 Data every Tclock

So instead of acquiring Data @ The output of the counter we can acquire it at the output of the comparator \*

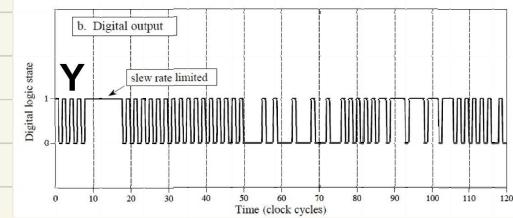
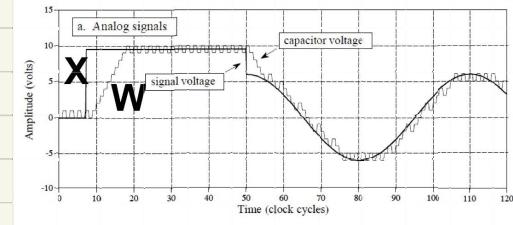
In this way this tracking ADC can be seen as 1bit oversampled ADC



What does the Δ modulator perform on the original signal?



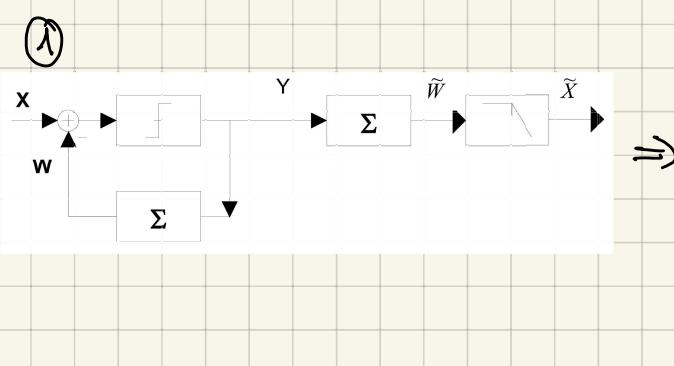
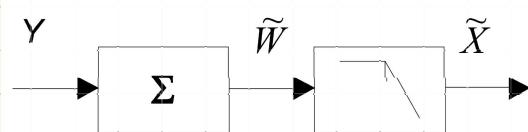
↑ it's an ADDER that behaves as a COUNTER which is incrementing or decrementing @ every clock period and then followed by a DAC



Notice: The D modulator can be considered as a 1st example of oversampling indeed we are not operating @  $f_s = 2f_{max}$  but @  $f_{clock} \gg f_s = 2f_{max}$

Then we can take the output stream  $\tilde{Y}$  and provide it to a counter which is identical to the original one and finally we need a LPF to reconstruct the original signal

### Delta Reconstructor

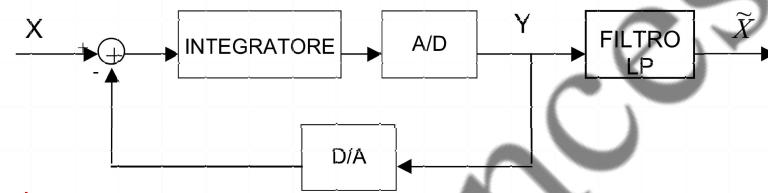


This configuration represents a  $\Sigma\Delta$  MODULATOR

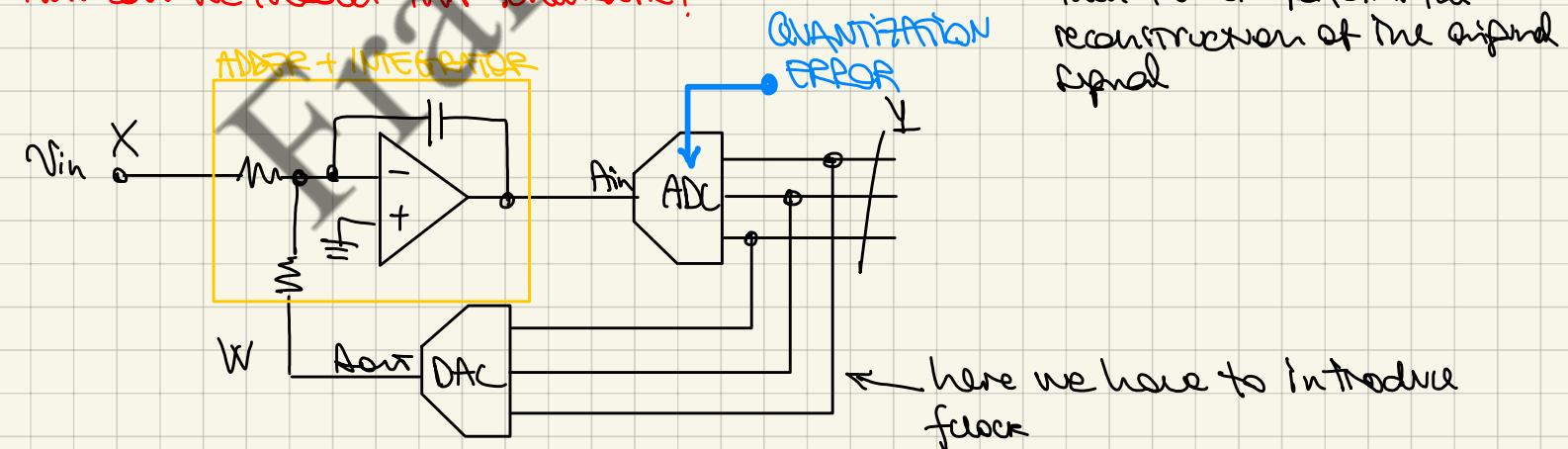
bcz the two  $\Sigma$  converge to the same input differentiator node. They can be placed in the same point w in the loop.

**SIGMA-DELTA MODULATOR** → it's an integrating network

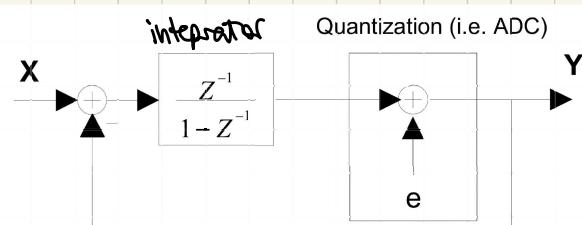
\* is no more the counter that we had at the beginning, but now it has a brand-new meaning:



How can we model this schematic?



**LINAR MODELING:**

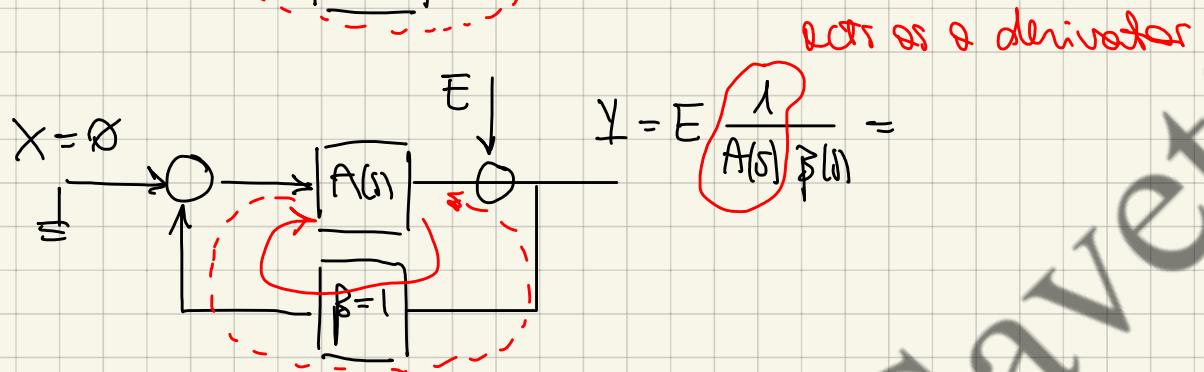
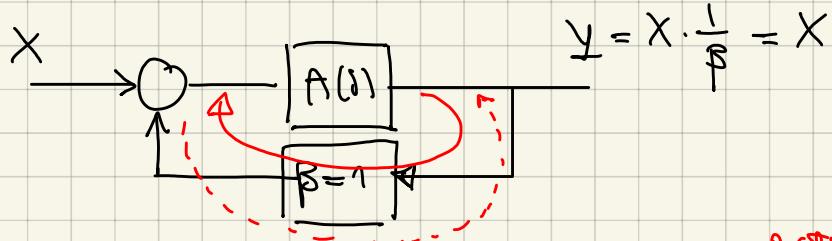
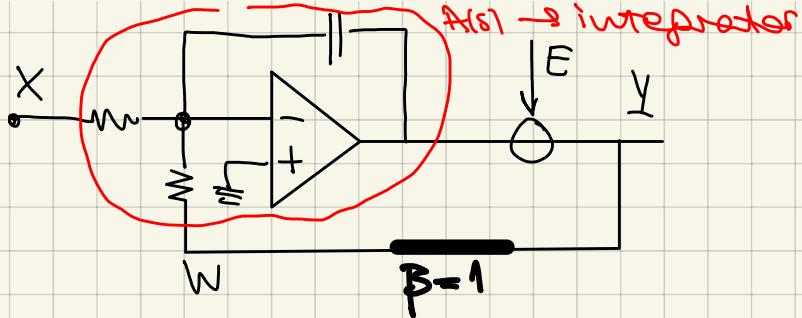


### GENERAL ARCHITECTURE

The input signal is processed by a quantizer after the integration. The quantized output  $Y$  is reported to the input by means of a FB and subtracted from the input signal.

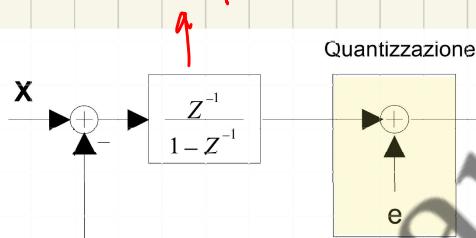
Then the LPF performs the reconstruction of the original signal

here we have to introduce clock



The loop contained in the circuit works as a modulator of the quantization noise but it reduces its presence in the band, shaping the spectrum to the higher freq.s  
 ↳ this effect is known as NOISE SHAPING

$Z$ -transform



$$Y(z) = X(z) \frac{1}{1 + \frac{1 - z^{-1}}{z^{-1}}} + e(z) \frac{1 - z^{-1}}{1 + \frac{1 - z^{-1}}{z^{-1}}} =$$

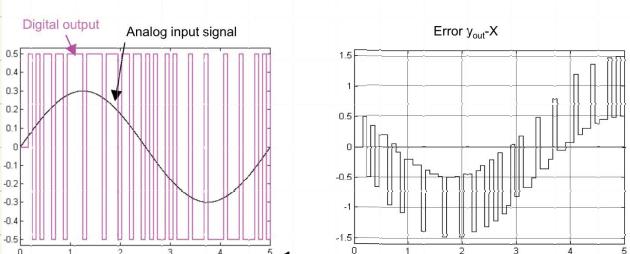
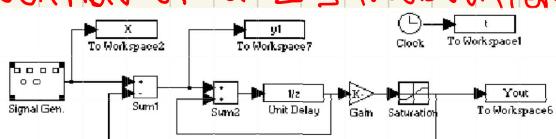
$$= \underbrace{X(z) z^{-1}}_{\text{Eq is still present but it gets differentiated}} + \underbrace{e(z) (1 - z^{-1})}_{\text{1 clock pulse delayed}}$$

Eq is still present but it gets differentiated

1 clock pulse delayed

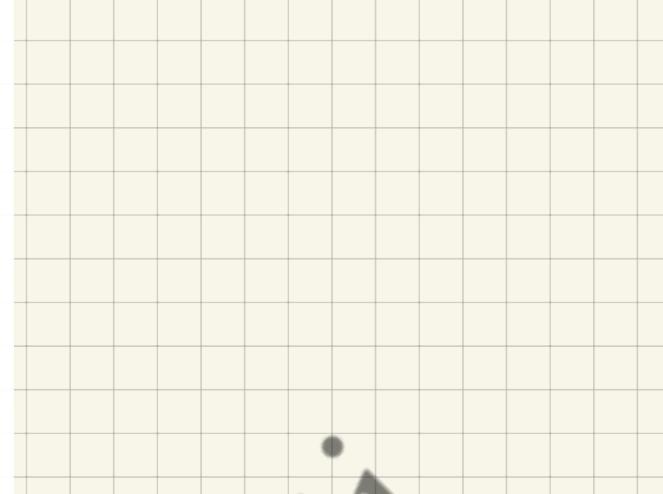
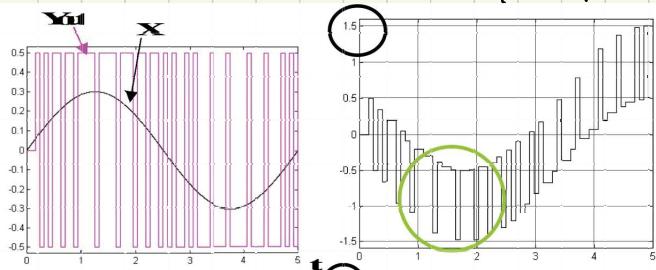
The noise shaping acts as  $(\sinh(1))^2$ :  $Q(f) = Q_0 |1 - e^{j2\pi f T}|^2 = 4Q_0 \sin^2(\pi f T)$

SIMULATION OF A  $\Sigma\Delta$  MODULATOR

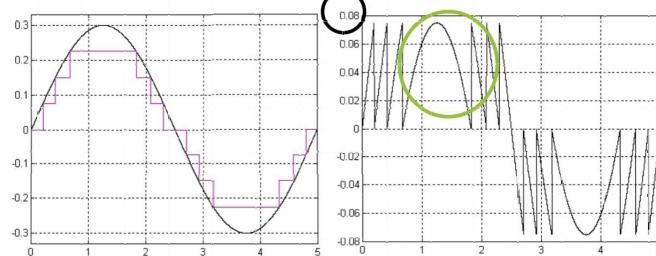


The error of The  $\Sigma\Delta$  modulator is very high but it stays @ very high frequency

Sigma-Delta:  
1bit



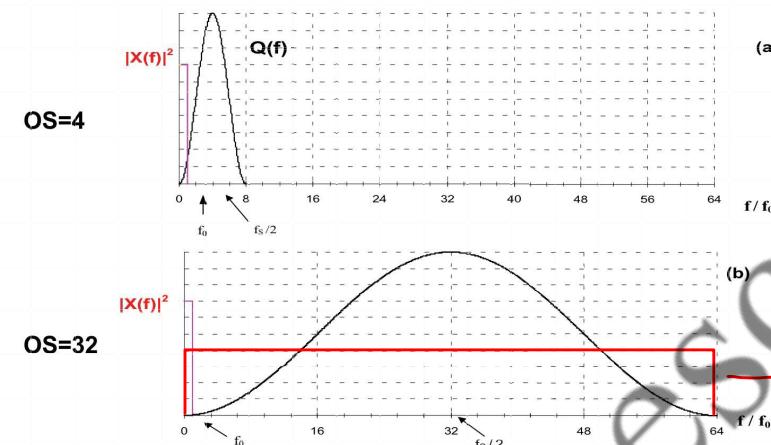
Classic ADC:  
3bit



Why we may say that The  $\Sigma\Delta$  modulator is eventually better than The classical oversampled ADC?

$$Y(z) = X(z) \cdot z^{-1} + e(z) \cdot (1 - z^{-1})$$

$$Q(f) = Q_0 |1 - e^{j2\pi fT}|^2 = 4Q_0 \sin^2(\pi fT)$$



→ Due to this shaping here, the quantization error's spectrum is no longer flat as it was in the standard oversampling technique.

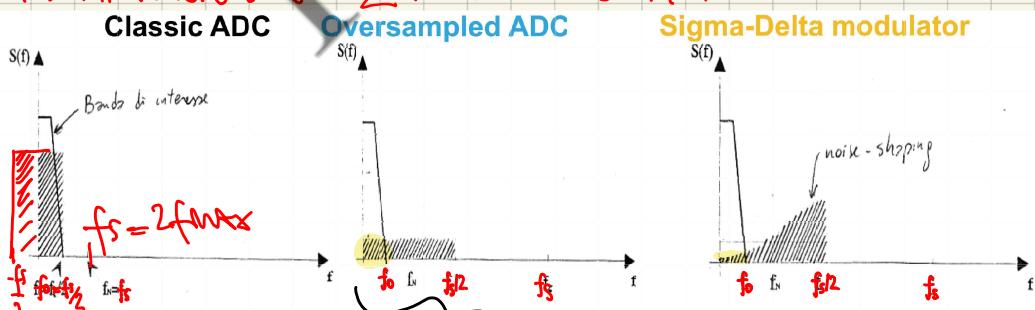
Thanks to The  $\Sigma\Delta$  architecture, and so to the fact that first there's an integration and then there'll be the ADC which introduce an error at the latter stage we have an excellent Noise Shaping

→ Eg if no longer flat bw f\_0 and f\_s happens w/ the standard oversampling

from the pov of the quantization noise, The  $\Sigma\Delta$  modulator improve the performance levels obtainable by the classical oversampling, thanks to the out-of-band noise shaping.

The noise is translated towards the high frequencies, increasing the oversampling factor.

## ADVANTAGES OF $\Sigma\Delta$ MODULATION



In the oversampled ADC we use a higher  $f_s$ , but we don't filter @  $f_s/2$ , we are smart and we filter @  $f_{max}$  ( $f_0$ )  
In this way we "acquire" only a portion of the total Eq

→ In fact, in this way we gain 1/2 bit every doubling of  $f_s$

In the Sigma-Delta modulator we get an even better result.

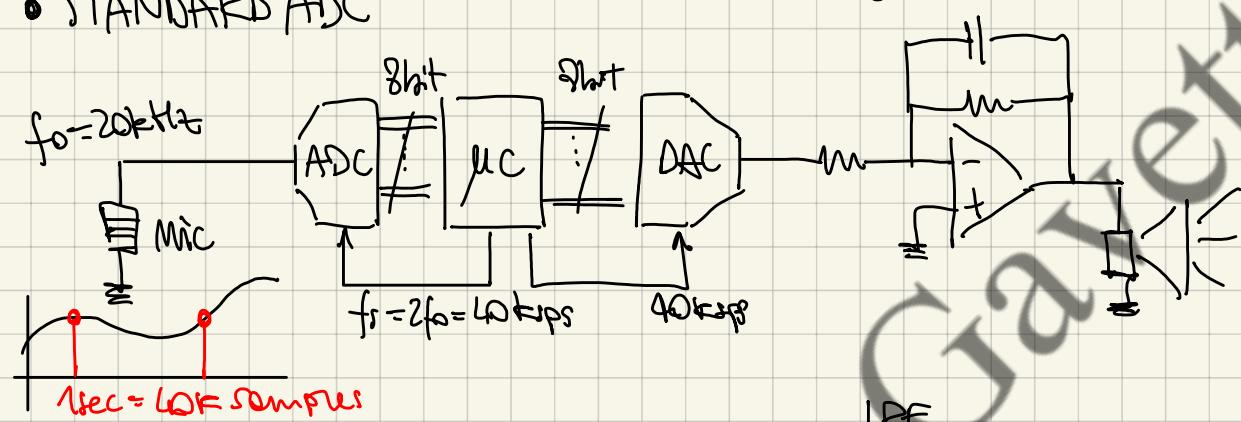
We sample @ an higher  $f_s$  w/r/t the standard ADC, as we do for the oversampled ADC we do not filter @  $f_s/2$ , but we are smart and we filter @  $f_0$ .

Here the advantage is that  $\text{Eq}$  gets a noise shaping so the portion of  $\text{Eq}$  we acquire is even lower w/r/t the amount we acquire w/ a standard oversampled ADC

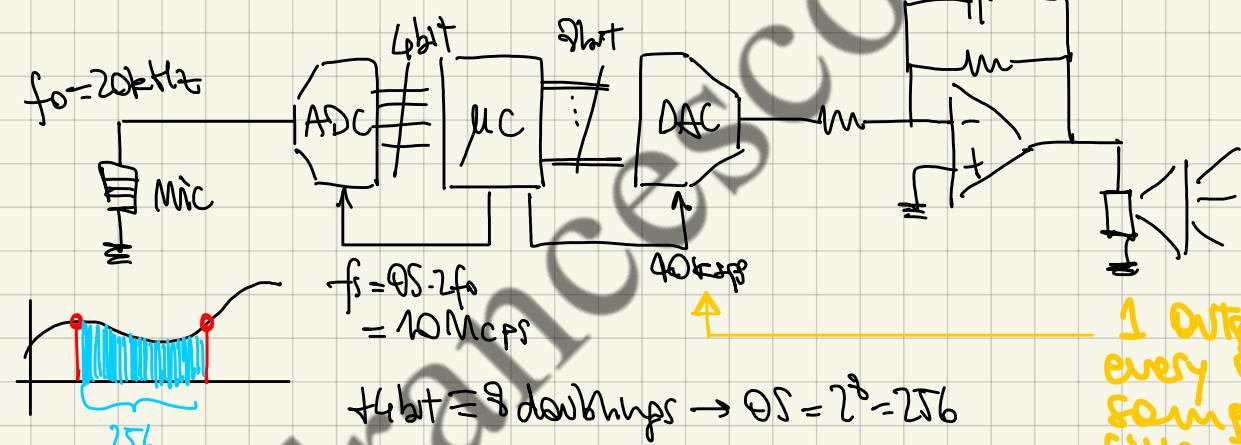
**RESULT:** Improvement of  $2^3 = 8 \text{ dB} = 1.5 \text{ bit}$  every doubling of  $f_s$  (i.e. doubling of OS)

EXAMPLE:

• STANDARD ADC



• OVERSAMPLED ADC



1 output sample every OS input samples and filter them w/ an average comparator

→ we have to use a larger memory!!!

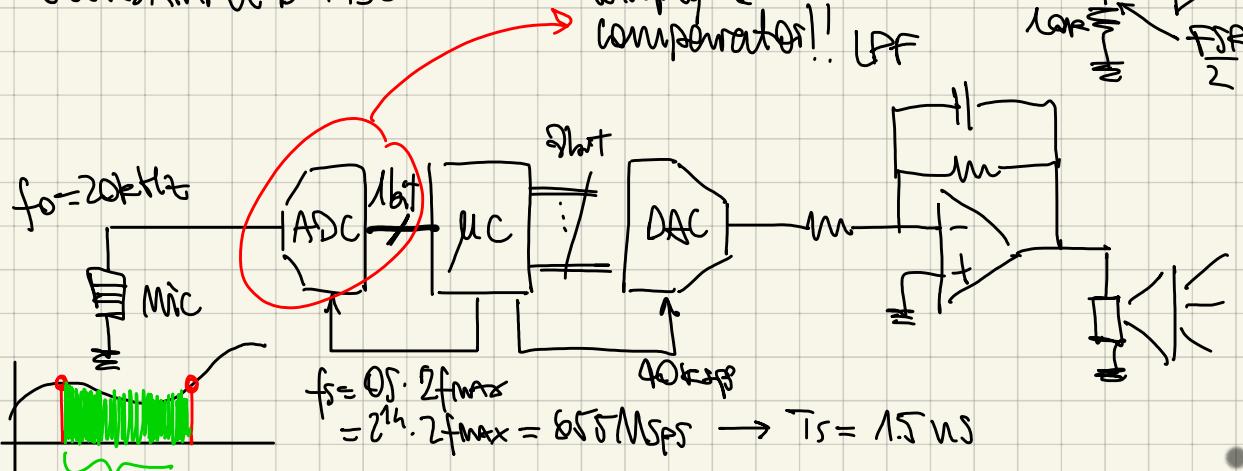
↳ 10 MB

16 16 16 16  
15 15 15 15  
14 14 14 14  
13 13 13 13  
12 12 12 12  
11 11 11 11  
10 10 10 10  
9 9 9 9  
8 8 8 8  
7 7 7 7  
6 6 6 6  
5 5 5 5  
4 4 4 4  
3 3 3 3  
2 2 2 2  
1 1 1 1  
0 0 0 0

result of the averaging

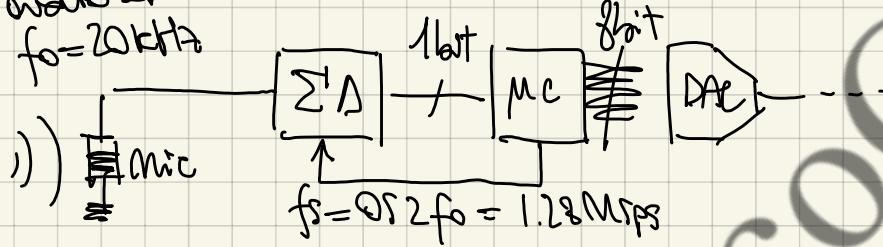
14.3

## • EXTREME OVERSAMPLED ADC



## • $\Sigma\Delta$ MODULATOR

oversampled



$$+7\text{bit} \equiv 5\text{ doublings} \rightarrow QOS = 2^5 = 32$$

$(+1.5\text{ bit every doubling}) \rightarrow +7.5\text{ bit actually}$

## $\Sigma\Delta$ MODULATOR'S RESOLUTION

$$\text{SNR}_{\text{dB}} = 6.02(c + 1.5L) - 3.41 \quad \text{where } c = \text{bit of } \Sigma\Delta$$

$$QOS = 2^L$$

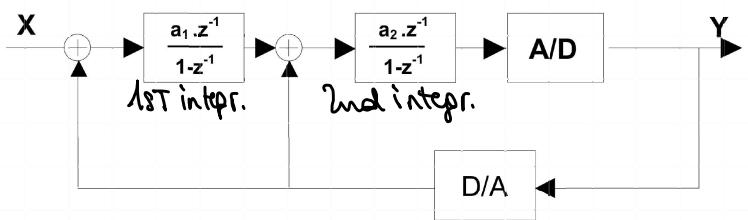
## MAXIMUM EQUIVALENT # OF BITS

$$b = \frac{\text{SNR} + 3.41}{6.02}$$

Any filtering/processing that provides a number of bits higher than b is useless, since the extra-bits are redundant, since they simply describe the quantization error!

## 2nd ORDER $\Sigma\Delta$ MODULATOR

Second-Order  $\Sigma\Delta$  modulators:

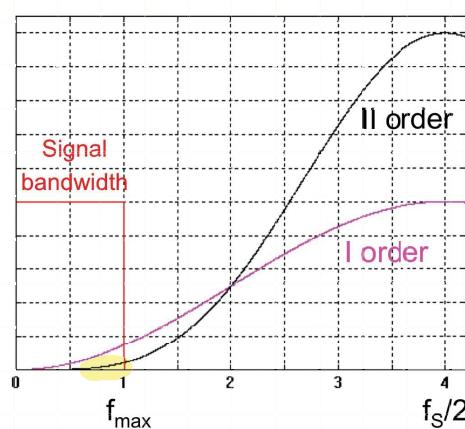


Notice: if the ADC is just 1 bit, so if it is a comparator, we don't need a DAC

$$Y(z) = X(z) \cdot z^{-2} + e(z) \cdot (1 - z^{-1})^2$$

→ the noise shaping acts no more as  $\sin^2(\theta)$  but as  $\sin^4(\theta)$ :

$$Q(f) = 16 Q_0 \sin^4(\pi f T)$$



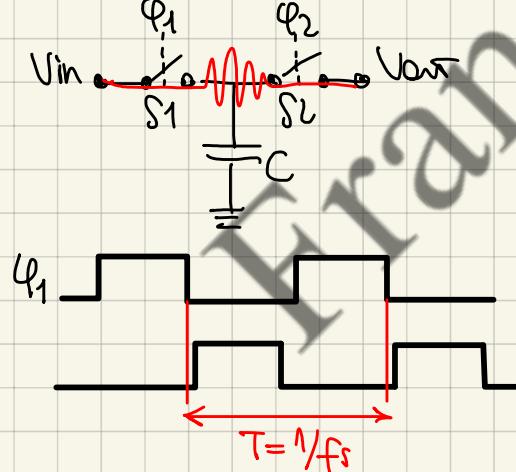
Filtering @  $f_0$ , we collect even bits Eq w/ respect to  
The 1st order  $\Sigma\Delta$  modulator

RESULT:

Improvement of  $Z = ?$  dB = 2.5 bit every doubling of  $f_s$   
(i.e. doubling of OS)

How is a  $\Sigma\Delta$  modulator designed?

- HOW TO DESIGN A R IN A TIME-DISCRETE WORLD:



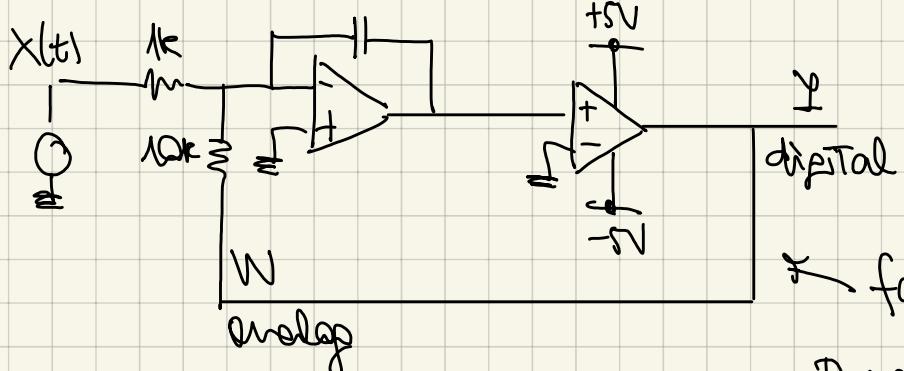
During each clock cycle we charge C to  $V_{in}$  and then we discharge it to  $V_{out}$ , so the charge transfer is:

$$Q = C(V_{in} - V_{out})$$

$$\rightarrow I = \frac{Q}{T} = f_s \cdot C (V_{in} - V_{out})$$

In the continuous-time operation:  $I = \frac{(V_{in} - V_{out})}{R}$

$$\Rightarrow R = \frac{1}{f_s \cdot C}$$



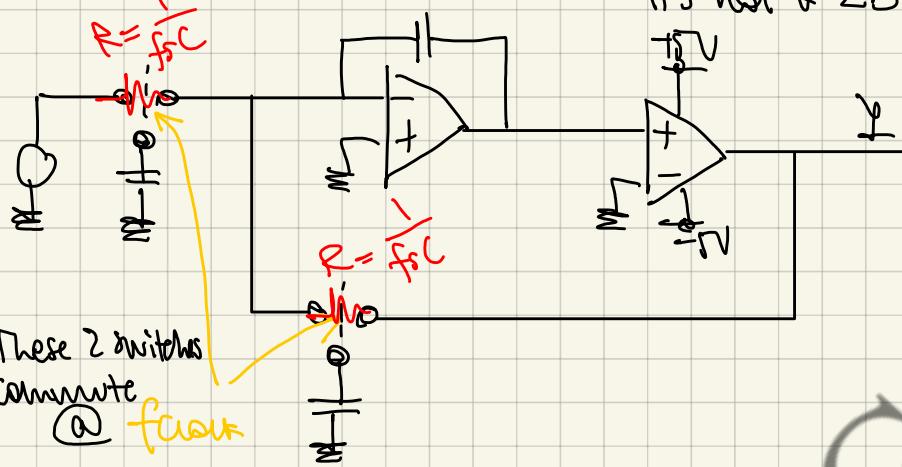
↙ flock?

There's ↘ no flock

it's not a ΣΔ modulator, it's an amplifier



⇒ This is a ΣΔ Modulator

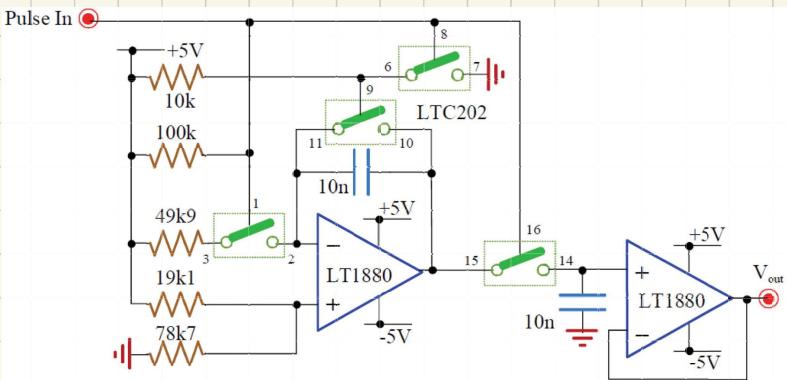


These 2 switches  
commute  
at flock

## ES16 - EXERCISES

15/12/2021

### EX3



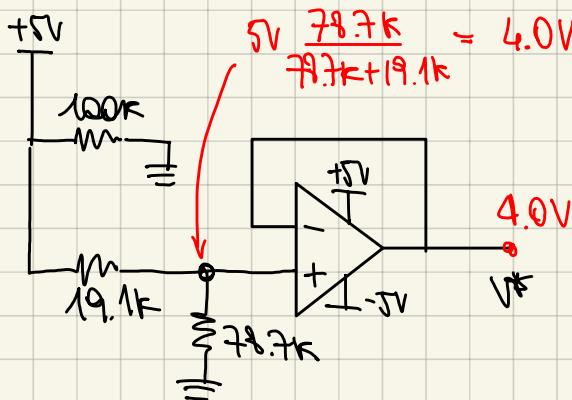
The LTC202 is a quad analog switch (closed when control pin is high)

The input is a pulse, whose width  $T_{pulse}$  is in  $1\text{ms} \div 2\text{ms}$  range.

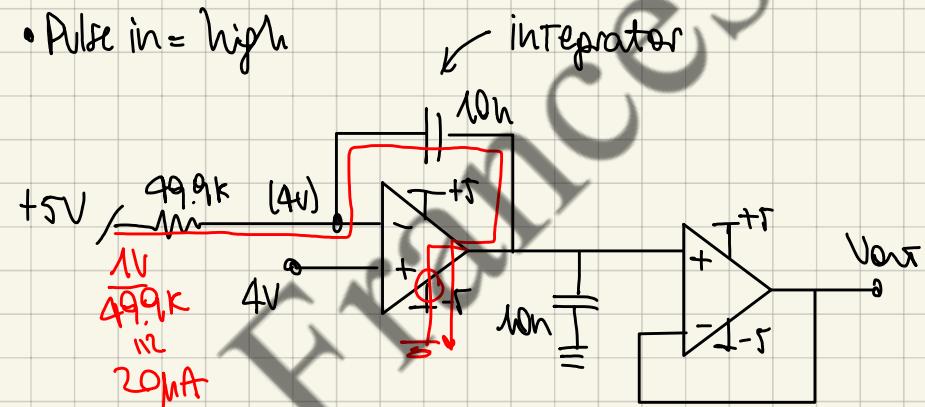
The opamps have:  $I_S < 1.5\text{nA}$  in the  $-40^\circ \div +75^\circ\text{C}$  range

(a) Compute  $V_{out}$  as a function of  $T_{pulse}$

• Pulse in = low

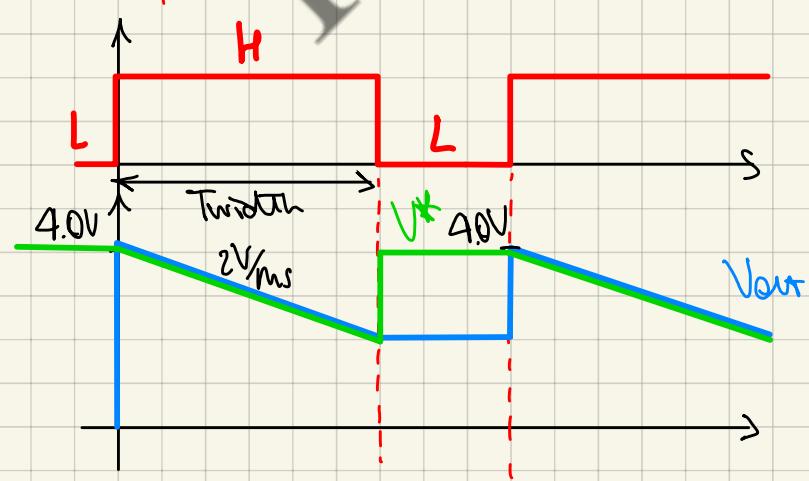


• Pulse in = high



$$\frac{dV}{dt} = \frac{i}{C} = \frac{20\mu\text{A}}{10\text{nF}} = 2\text{kV/s}$$

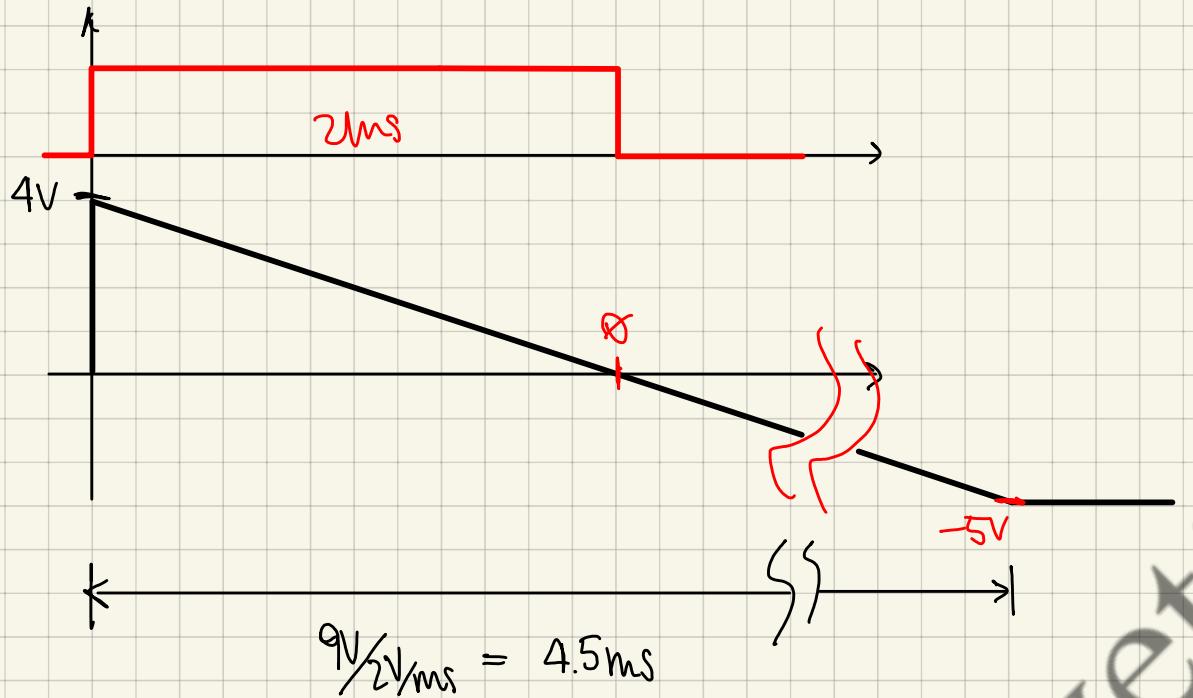
$$= 2\text{V/ms}$$



$$V_{out}(T_w) = 4V - \frac{2V}{ms} \cdot T_w$$

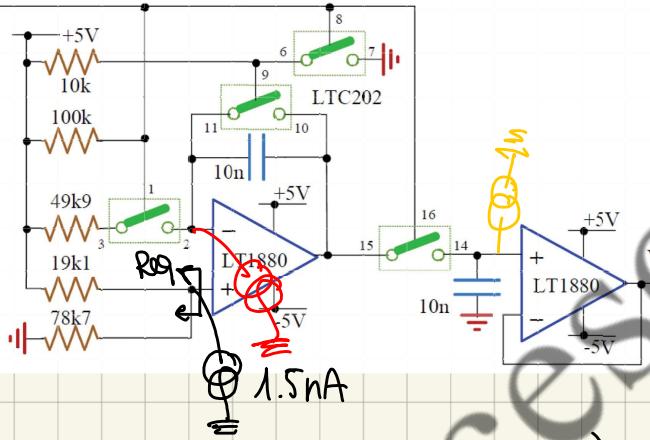
The output gives info about the duration of the pulse

There's a limit



③ Reckon the min. time that guarantees a precision of 1μs.

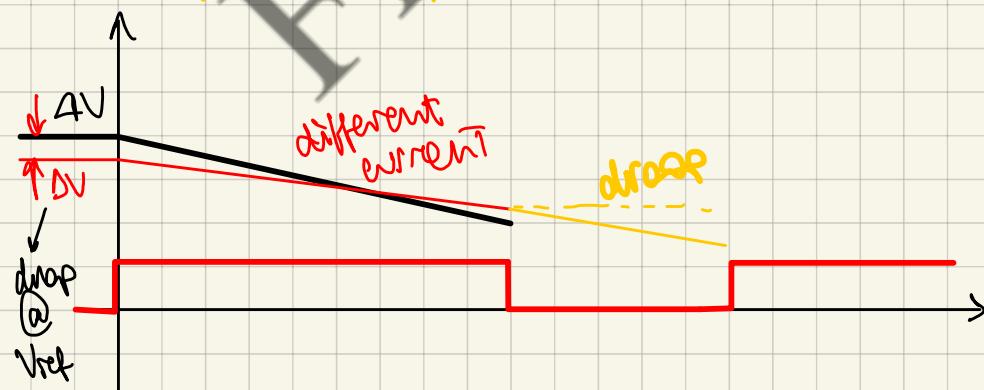
Pulse In



$$\Delta V = 1.5 \text{nA} \cdot R_B = 1.5 \text{nA} \cdot (19.1 \text{k} \parallel 78.7 \text{k}) = 15 \text{nA} \cdot 15.1 \text{k} = 23 \mu\text{V}$$

$$i = 20 \mu\text{A} - I_B =$$

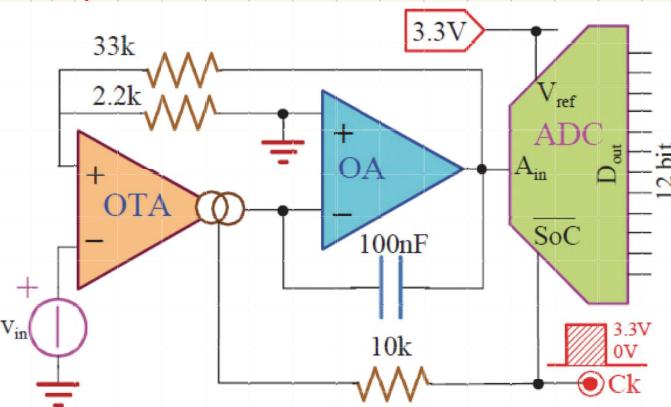
$$\text{droop} = \frac{dV}{dt} = \frac{I_B}{C} = \frac{1.5 \text{nA}}{10 \text{nF}} = 0.15 \frac{\text{V}}{\text{s}}$$



## ES16 - EXERCISES (3)

22/12/2021

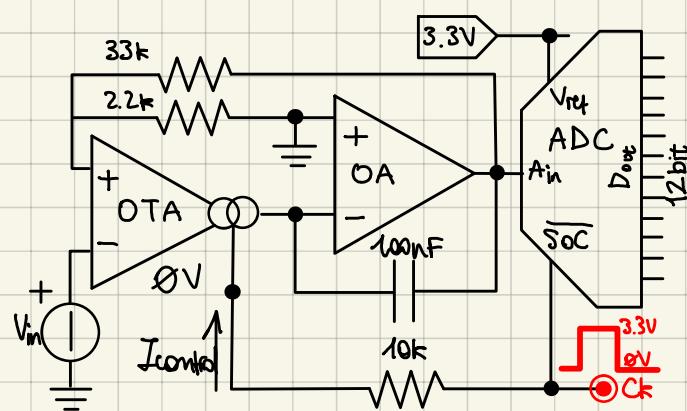
### EX12



12 bit, 3.3V FSR ADC with  $1\text{mV}_{\text{rms}}$  total input noise (at its  $A_{in}$  input). Ideal OpAmp. OTA control pin at 0V.

- For a 5mV peak sinusoidal input with 5mV baseline, compute the  $\text{SNR}_{\text{ideal}}$ ,  $\text{SNR}_{\text{theor}}$ ,  $\text{SNR}_{\text{real}}$ , when considering just the  $1\text{mV}_{\text{rms}}$  noise at the ADC input.
- Compute the acquisition bandwidth and  $t_{\text{acquisition}}$  for maximum FSR and  $\frac{1}{2}\text{LSB}$  error.
- Quote the static error in LSB due to  $I_B=1\mu\text{A}$  and  $V_{os}=2\text{mV}$  of the OpAmp.

(A)



$$\bullet C_{lk} = \emptyset$$

when  $C_{lk} = \emptyset$ , since control pin  $\bar{Q} = 0\text{V}$ ,  $I_{\text{control}} = 0\text{A}$

$$G_m = \frac{I_{\text{control}}}{V_{th}} = \emptyset \rightarrow \text{OTA is off}$$

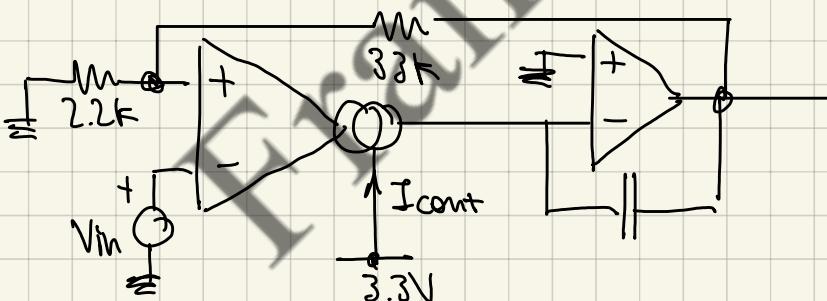
$$I_{\text{out}} = \emptyset$$

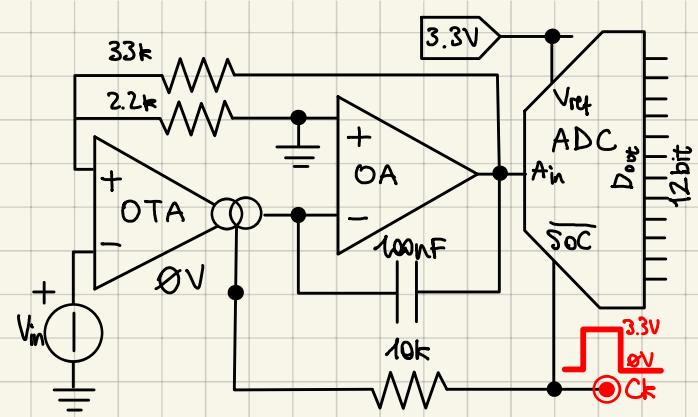
This means that the voltage stored into the capacitor stays constant and so the circuit acts as a  $\text{S/H}$  stage in HOLD phase

$\bullet C_{lk} = 3.3\text{V} \rightarrow \text{SAMPLING PHASE!}$

$$I_{\text{control}} = \frac{3.3\text{V}}{10\text{k}} = 330\mu\text{A} \rightarrow$$

$$G_m = \frac{330\text{nA}}{25\text{mV} (\text{@ } R_1)} = 13.2 \text{ mA/V}$$





Francesco Gavetti