

„O viață liniștită și modestă aduce mai multă bucurie decât urmărirea succesului asociat cu anxietatea constantă”, Einstein a scris aceste cuvinte în germană pe antetul hotelului Imperial Tokyo.

ASRSV – curs 7

Sinteza Semnalului Vocal

- **Introducere**
- **Sisteme de sinteza clasice**
- **Tehnici de sinteza**

http://research.spa.aalto.fi/publications/theses/lemmetty_mst/chap2.html

Development of speech synthesizers

https://www.cs.cmu.edu/~srallaba/Learn_Synthesis/intro.html

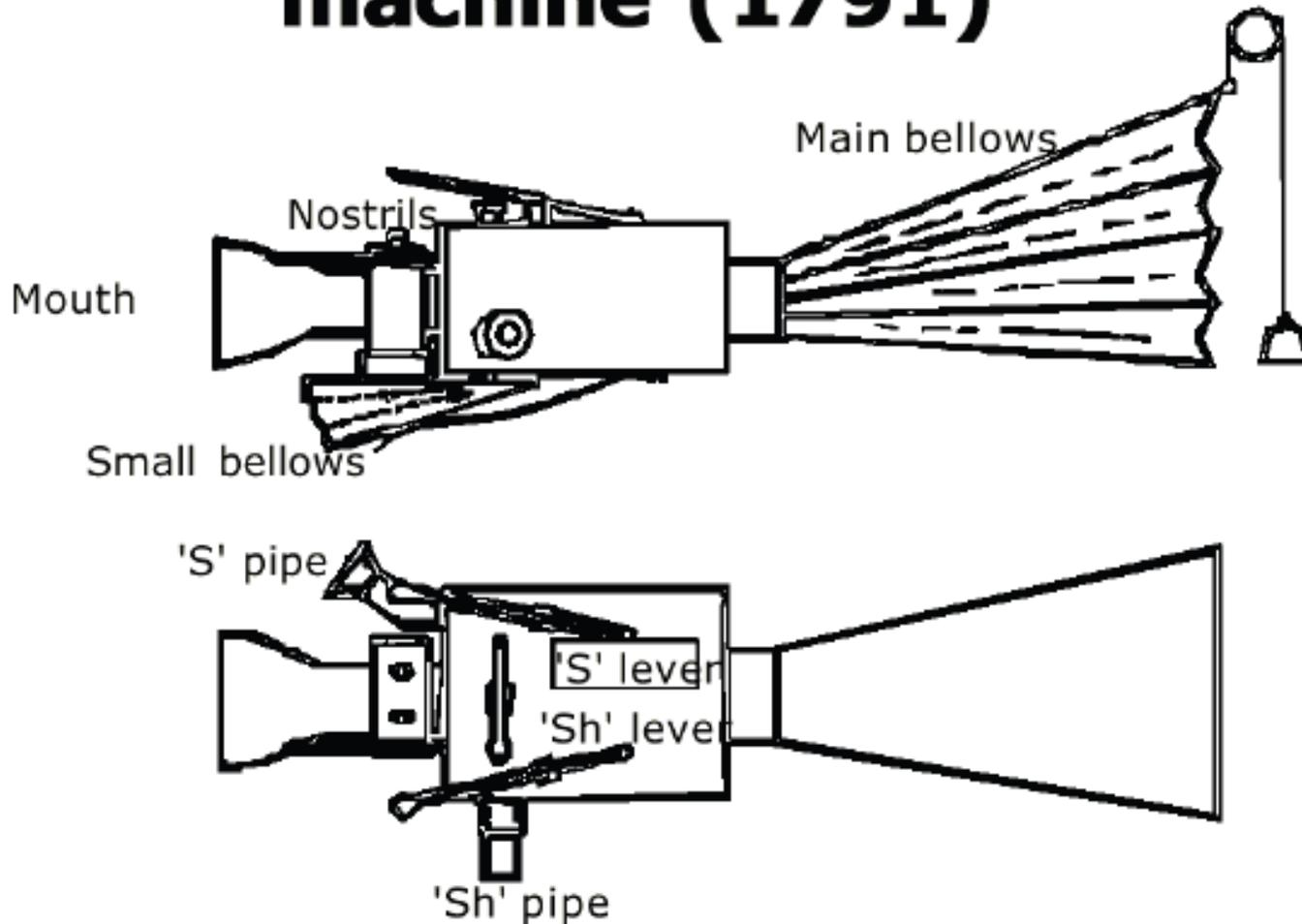
Sinteza vocală este generarea semnalului vocal sintetic utilizând intrare text (TTS) și modele care capturează informații fonetice, prosodice și acustice pentru a reproduce vorbirea umană într-un mod cât mai natural și inteligibil.

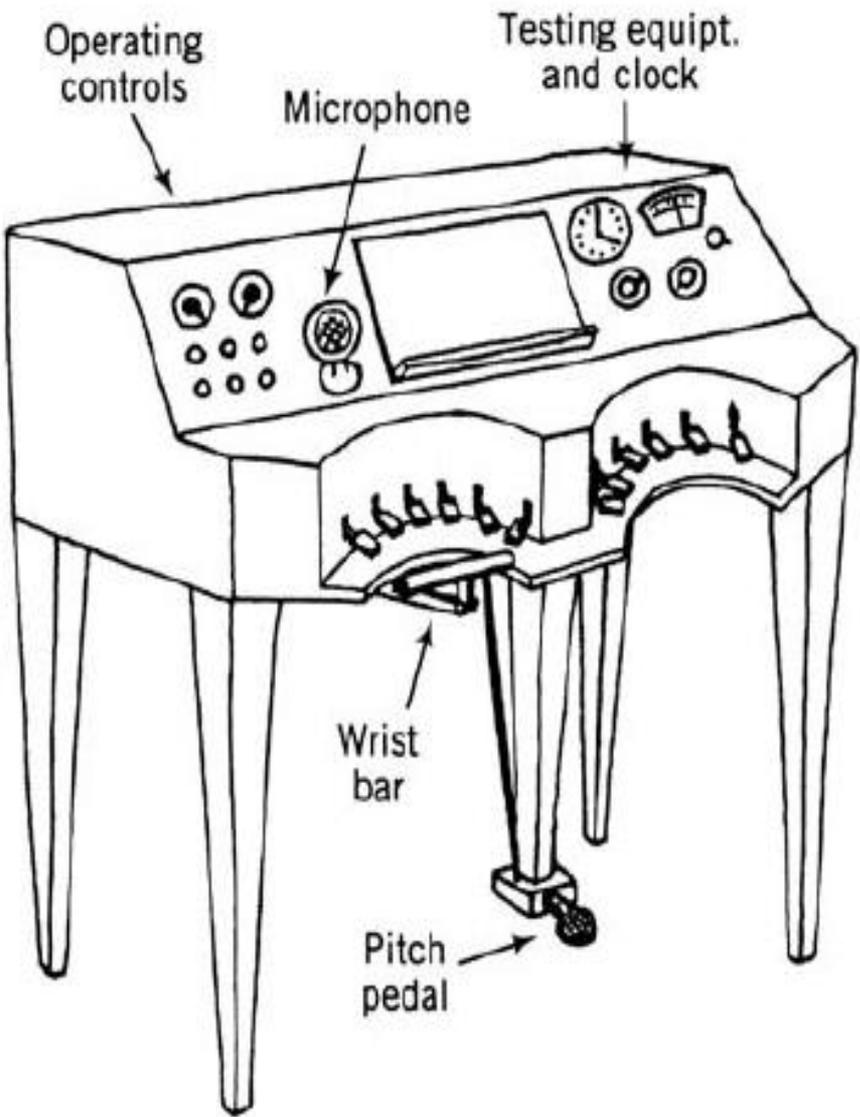
Obiective tehnice: se urmărește atât fidelitatea perceptuală (naturalitate, timbru, prosodie), cât și robustețea în diferite limbi, accente și voci, cu resurse de calcul compatibile cu aplicațiile în timp real.

Arhitecturi majore:

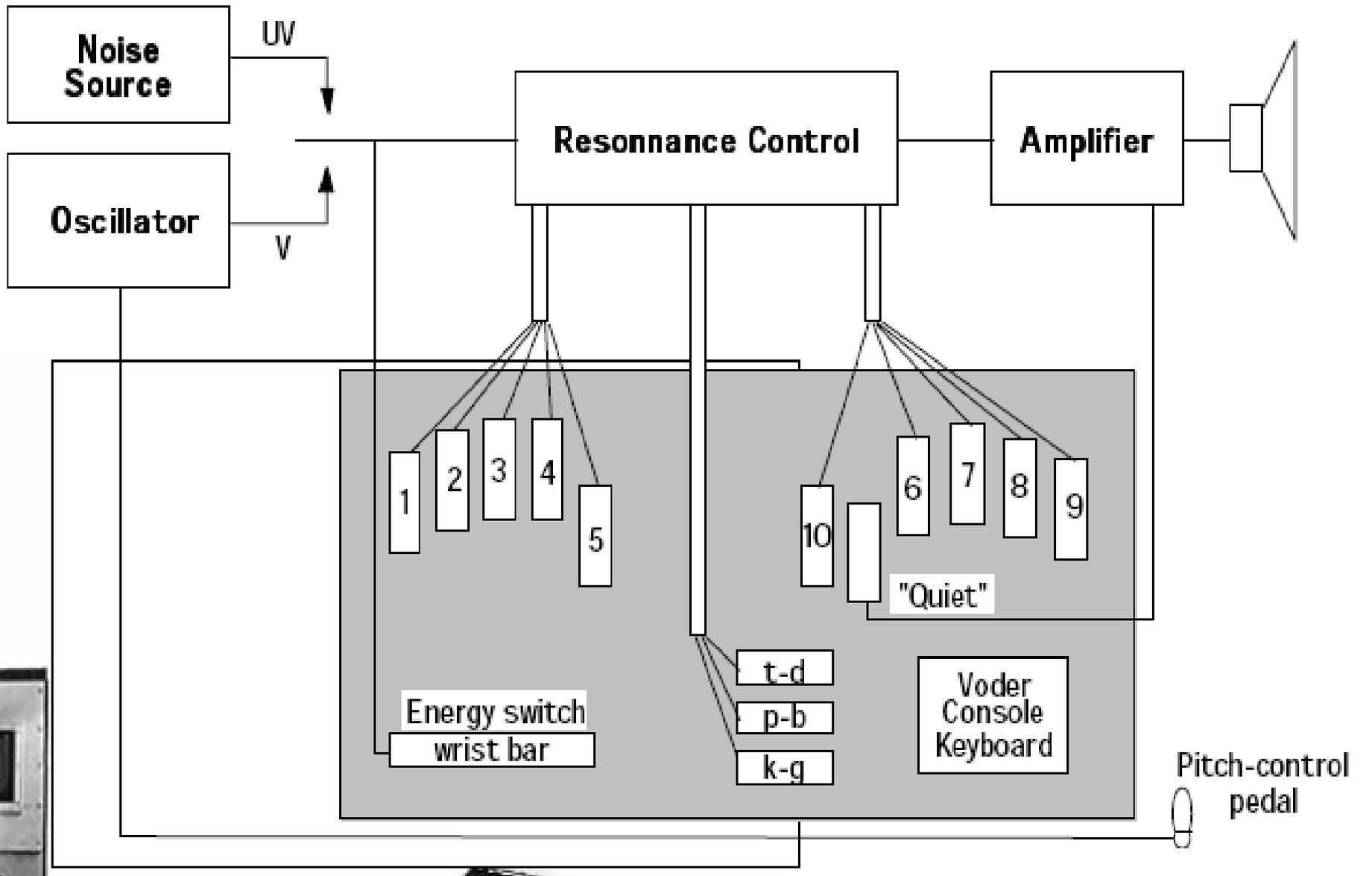
- pre-electronic (sinteza mecanică),
- electronică (formant-based), concatenativă (unit selection), parametrică și cele neuronale end-to-end actuale (Tacotron, WaveNet, transformere audio).

Von Kempelen's talking machine (1791)

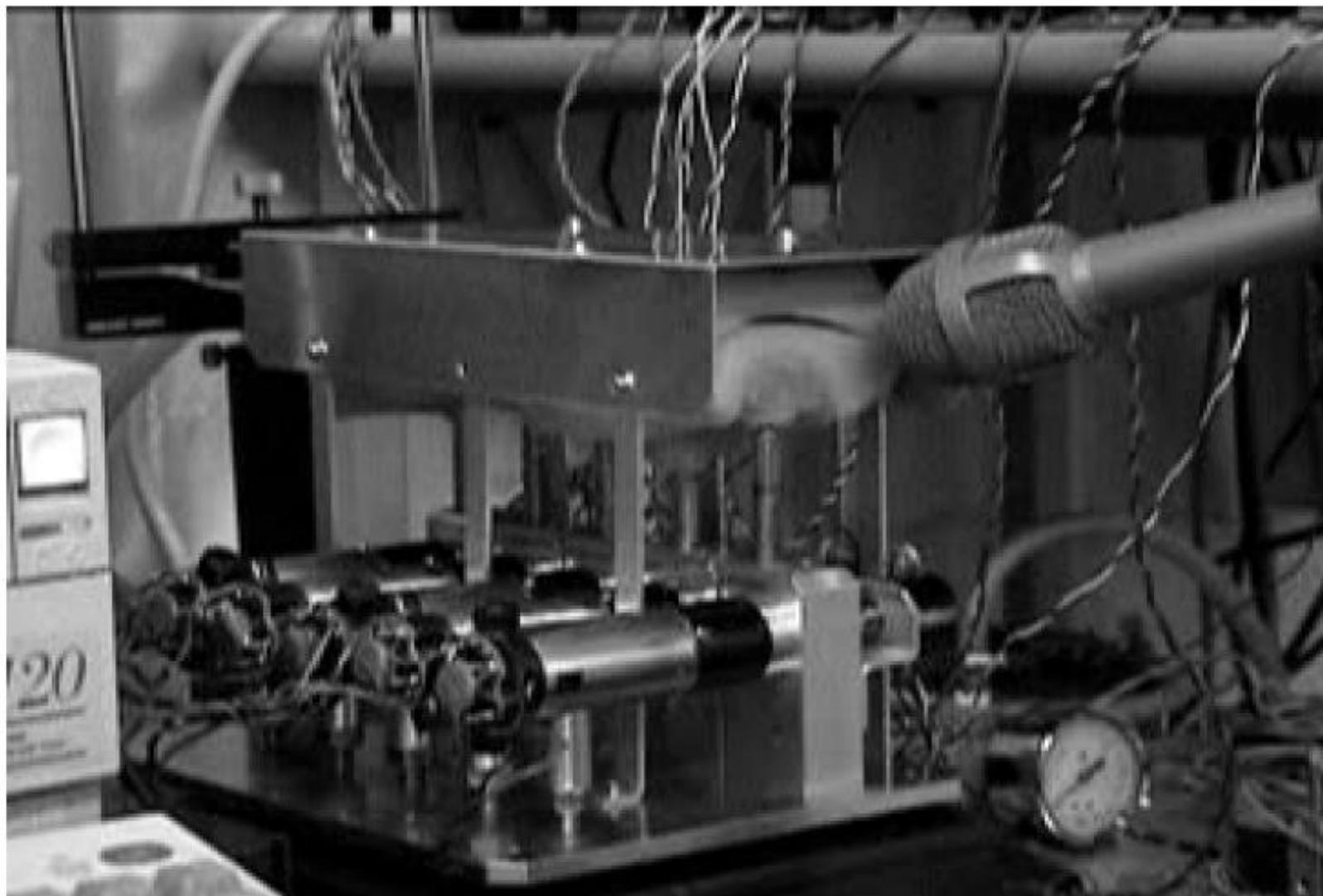




1939 –VODER- Homer Dudley (Bell Labs.)
(Voice Operated DEMonstrator)



Cercetari in sinteza electro-mecanica a vorbirii



<http://www.eng.kagawa-u.ac.jp/~sawada/>

Tehnologia sintezei vorbirii bazată pe învățare profundă

- Odată cu dezvoltarea științei și tehnologiei informatice, inteligibilitatea și naturalețea vorbirii sintetizate au fost mult îmbunătățite datorită progresului continuu al tehnicilor de sinteza TTS:
 - de la metodele bazate pe formanți
 - la metodele de concatenare a formelor de undă
 - metode bazate pe selecția unităților de sinteza
 - la metodele de sinteză statistică parametrică a vorbirii (SPSS) bazate pe modele Markov ascunse (HMM)
- La metodele bazate pe ‘deep-learning’ care :
 - [Learning from Audio: Wave Forms](#)
 - [Learning from Audio: Time Domain Features](#)
 - [Learning from Audio: Fourier Transformation](#)
 - [Learning from Audio: Spectrograms](#)
 - [Learning from Audio: Pitch and Chromagrams](#)

THE BEST TEXT-TO-SPEECH SOFTWARE

Click the links below to go to the provider's website:

1. Amazon Polly
2. Linguatec Voice Reader
3. Capti Personal
4. NaturalReader
5. Voice Dream Reader

Or, jump to:

[Best free text-to-speech apps.](#)

- [1. Amazon Polly](#)
- [2. Linguatec Voice Reader](#)
- [3. Capti Personal](#)
- [4. NaturalReader](#)
- [5. Voice Dream Reader](#)

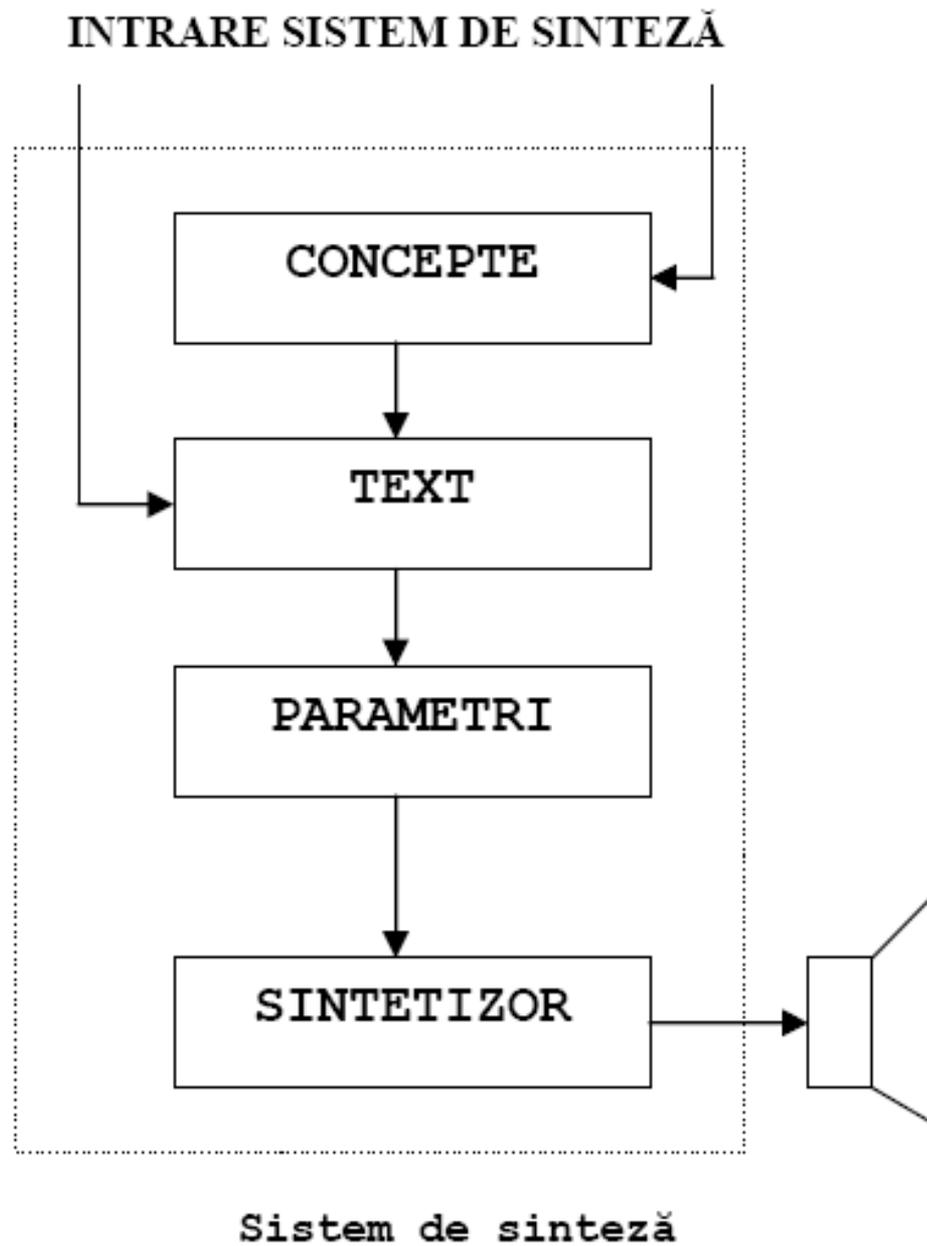
jump to:

[Best free text-to-speech apps.](#)



- **Alexa** - serviciul de voce cloud-based al companiei Amazon, disponibil pe zeci de milioane de dispozitive de la producătorii de dispozitive Amazon și de la terți.
- Cu Alexa, putem experimenta un dialog natural care oferă clienților o modalitate mai intuitivă de a interacționa cu tehnologia pe care o folosesc în fiecare zi.

Sisteme de sinteza vocala clasice



Moduri de sinteza:

- **mod grosier** (*mem.mare, vocabular limitat*)
- **prin fragmentarea vorbirii** (*difoneme, demisilabe, reguli*)

Sintetizoarele se pot clasifica după *domeniul* în care se realizează *prelucrările* (timp sau frecvență)

- sintetizoare pe bază de eșantioane (amplitudine-timp)
- sintetizoare utilizând modulația delta
- sintetizoare utilizând generatoare de frecvență
- sintetizoare de canal (vocoder de canal)
- sintetizoare pe bază de predicție liniară (cu excitație F_0 , reziduală sau multiimpuls)
- sintetizoare formantice
- sintetizoare complexe (bazate pe forme de undă) cu eliminarea redundanței vorbirii

După **tehnica** folosită:

1 **generare directă** (codarea formei de undă) – include tehnici de prelucrare a formei de unda a SV înregistrat (direct/codat) => generarea mesajelor vocale

2 generarea vocii conform unui **model (analiză-sinteză)** - realizate prin modelarea producerii SV, la care vocea => parametri => sintetizorul => mesajul dorit. (sintetizoarele formantice, LPC,...)

3 **simularea tractului vocal**

- *metoda analogiei cu tractul vocal* care simulează propagarea undei acustice prin tractul vocal;

- *metoda analogiei terminale* care simulează structura spectrului cu caracteristicile de rezonanță și antirezonanță reproducând procesul articulator.

4 **generare pe bază de reguli** – SV => pe baza reg. fonetice, lingvistice pornind de la litere/grupuri fonemice + caracteristici prozodice

5 **metode neuronale end-to-end** actuale (Tacotron, WaveNet, transformere audio).

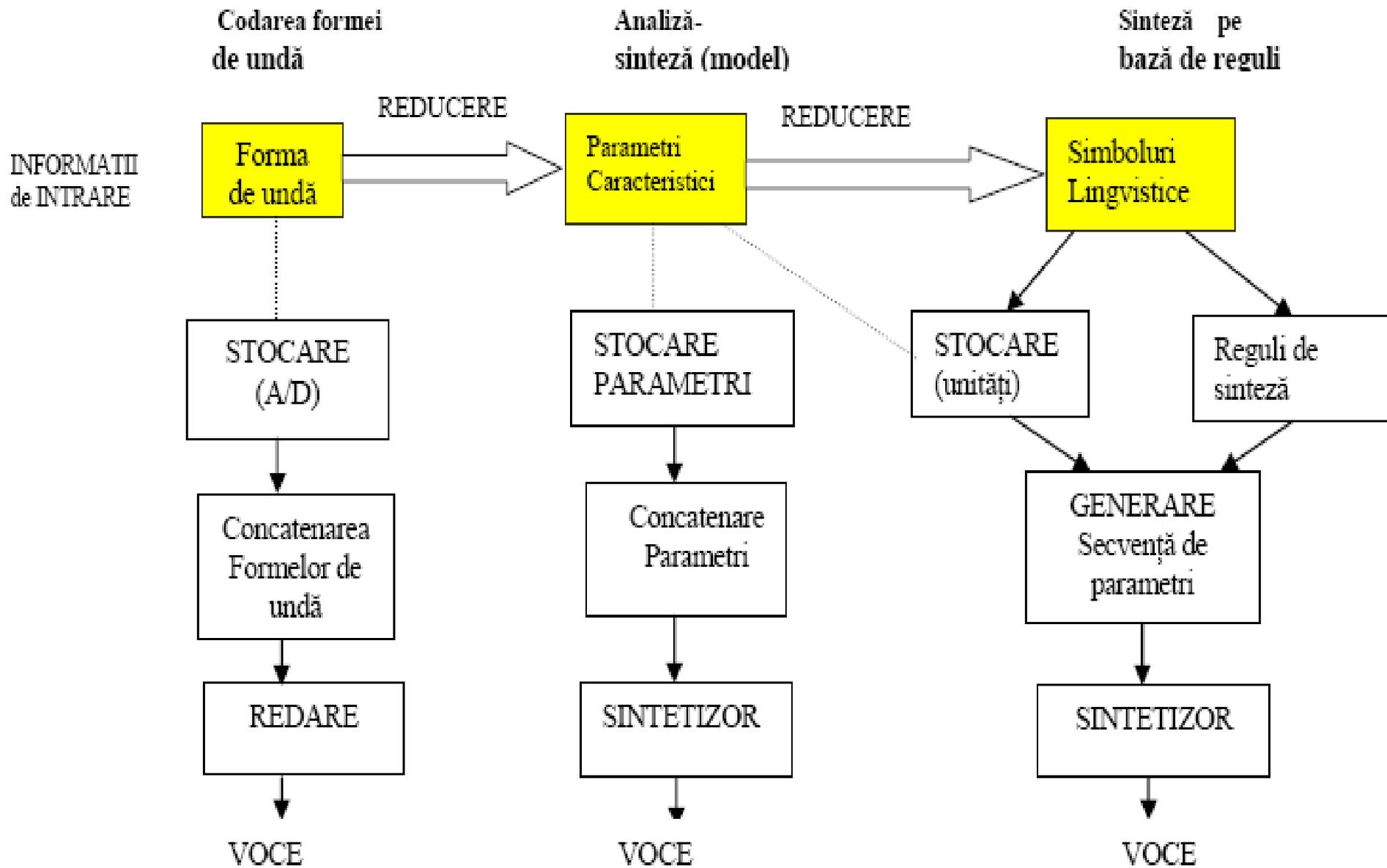


Fig2. Principiile de bază ale metodelor de sinteză și corelațiile între ele

CARACTERISTICI	CODAREA FORMEI	ANALIZĂ-SINTEZĂ	SINTEZĂ PRIN REGULI
INTELIGIBILITATE	Mare	Mare	Medie
NATURALETE	Mare	Medie	Medie
MARIME VOCABULAR	Mică(<500)	Mare >(1000)	Fără restricții
DEBIT BINAR	24-64Kbps	2.4-9.6Kbps	50-75bps
SEMNAL STOCAT ÎNTR-UN MBIT	15-40s	100s-7min.	Fără restricții
UNITĂȚI STOCATE	Cuvinte,silabe, propoziții	Cuvinte,silabe, propoziții	Foneme, silabe (CV,VCV,CVC, morfeme,..)
COMPLEXITATE	Redusă	Medie	Mare
RESURSE	Memorie	Procesor+Memorie	Procesor+Memorie
EXEMPLE TIPURI	PCM,ADPCM	Vocoder canal,LPC,..	Concatenare formă de undă, Vocoder (canal/cepstral/LPC)

Tabel 1. Caracteristici și performanțe ale tehnicilor de sinteză

Criteria de calitate pentru evaluarea sintezei

- **Inteligibilitatea**

(cuvinte, propoziții)

- **Naturațea**

(calitatea modelarii unităților ; prozodia, accentuarea, ezitățile etc)

Fluiditatea

(cursivitatea)

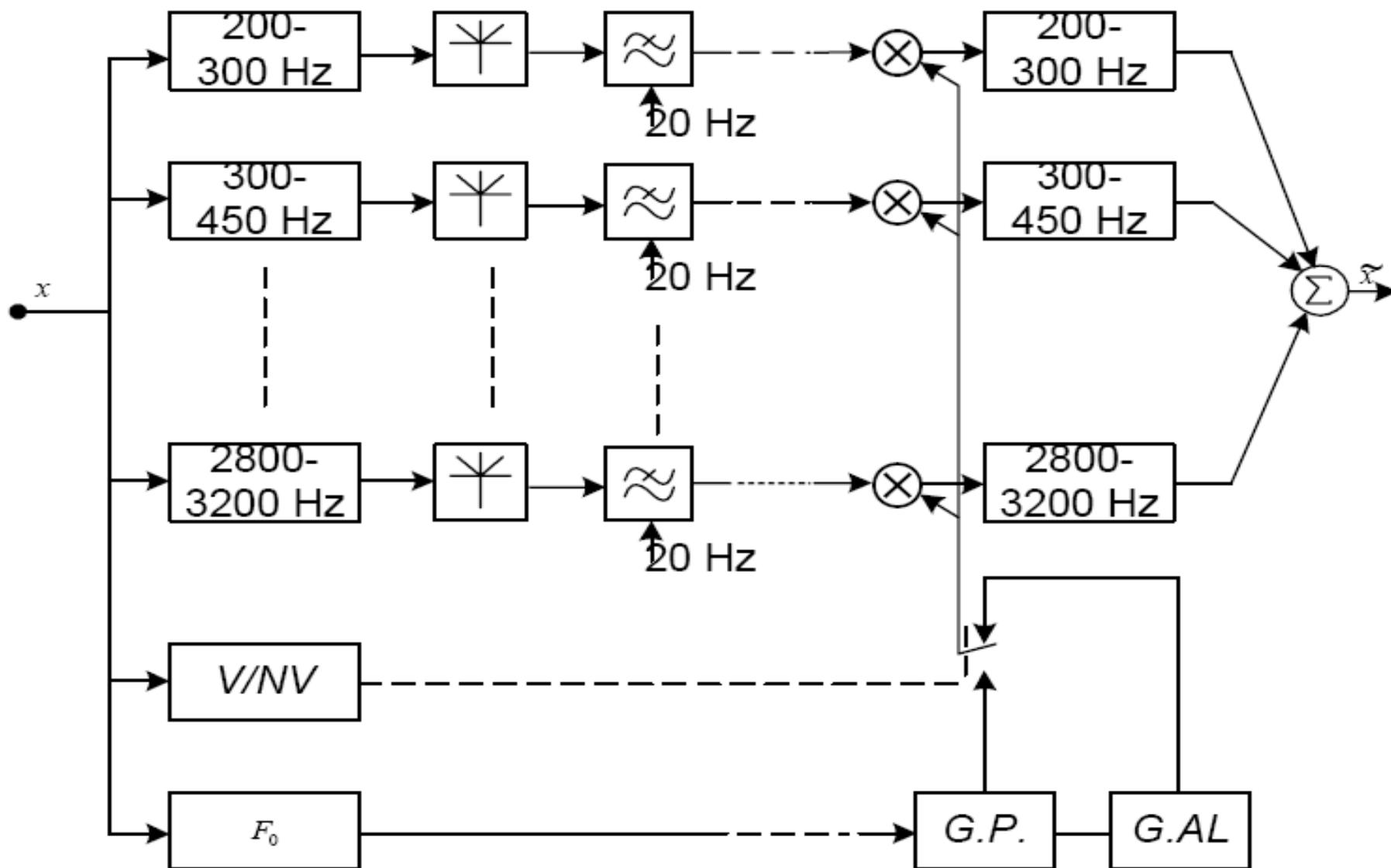
- **Potrivirea prozodiei**

(Expresii, continuări, întrebări, exclamatii)

- **Măsurători perceptuale:** MOS (mean opinion score), SIM (similaritatea percepută), preference tests; evaluări în diferite limbi pentru inteligibilitate și naturațea.
- **Măsurători obiective:** distanțe între spectre (MSE pe spectrogramă), dinamica F0, articularea (corelații cu seturi fonetice), statistici de intelligibilitate bazate pe n-gram –
- **Evaluare în timp real:** latențe, consum energetic, flexibilitate în personalizarea vocii, scalabilitate pe mai multe voci.
- **Criminalistică și fiabilitate:** detectabilitatea sintezei, metode de watermarking audio, reduceri ale riscurilor de impersonare/clonare.

1. SINTETIZOARE CU GENERARE DIRECTĂ (CODAREA FORMEI DE UNDĂ)

- Generează direct forme de undă sonoră pe baza unor măsurători efectuate asupra unor semnale reale



Sintetizor de canal

- Sinteza constă în concatenarea directă a unui număr oarecare de segmente SV
- Segmentele sunt anterior memorate, apoi sunt concatenate câteva dintre acestea, formând fonemele, pe urmă din acestea se formează silabe, morpheme ,... cuvintele.

Sintetizorul de canal – *consta în parametrizarea anvelopei spectrului SV pe termen scurt*

- se folosește un set de filtre (FTB 10-20) (sintetizor de canal / Channel Vocoder);

- **estimarea sonor/nesonor => CODATE**

- **F0**

- **RECEPTIE** => reconstituirea cu un set FTB \equiv , excitat cu F0/zg, redresat, FTJ

Banda $\sim 20\text{Hz}$ – anvelopa sp. => 300Hz

Avantaje:

- simplitatea și memoria relativ redusă necesară stocării datelor
- calitatea și inteligibilitatea este acceptabilă în anumite aplicații fără pretenții

Dezavantaje

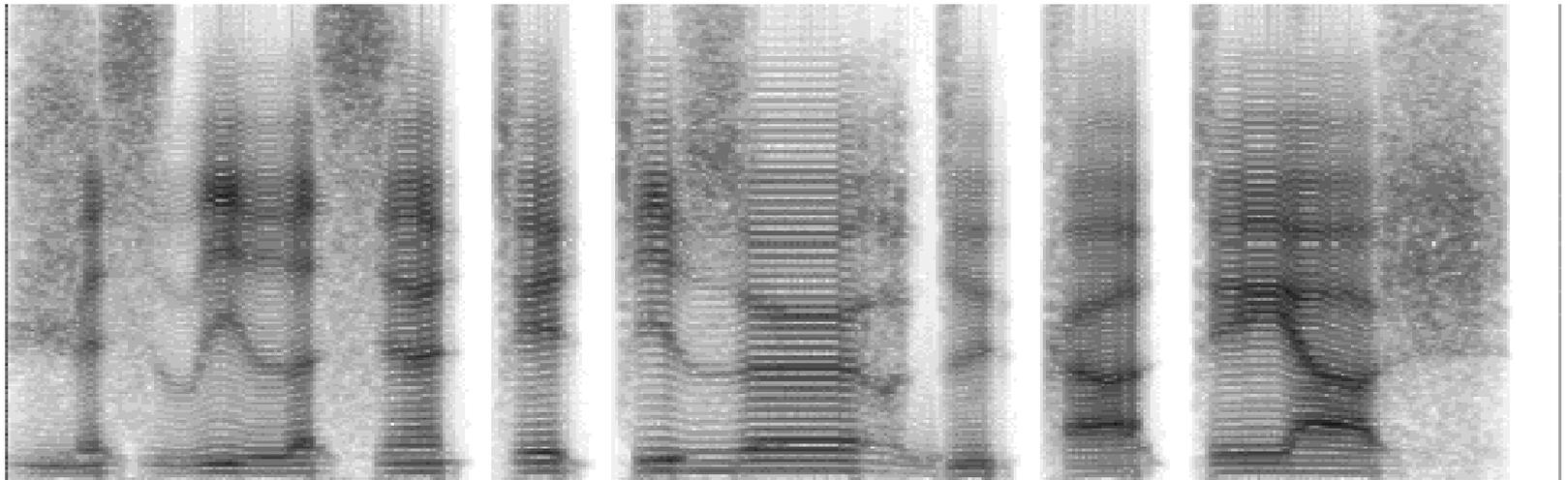
- descrierea spectrului este o aproximație mai mult sau mai puțin grosieră în funcție de numărul de canale utilizate, de lățimea lor de bandă

2. SISTEME DE SINTEZĂ A VOCII PE BAZA UNUI MODEL

SINTEZA FORMANTICĂ – 1950, Lawrence și Fant.

- În timp ce SV parcurge distanța de la glotă până la buze, spectrul larg al excitației este modulată datorită selectivității în frecvență a tractului vocal, rezultând formanții
- semnalul este generat printr-un sistem, *care simulează transmitanța tractului vocal, definită prin frecvențele de rezonanță (formanți-rez. Ord.2)*- caracterizat prin frecvența de rezonanță F_k și banda de trecere B_k .

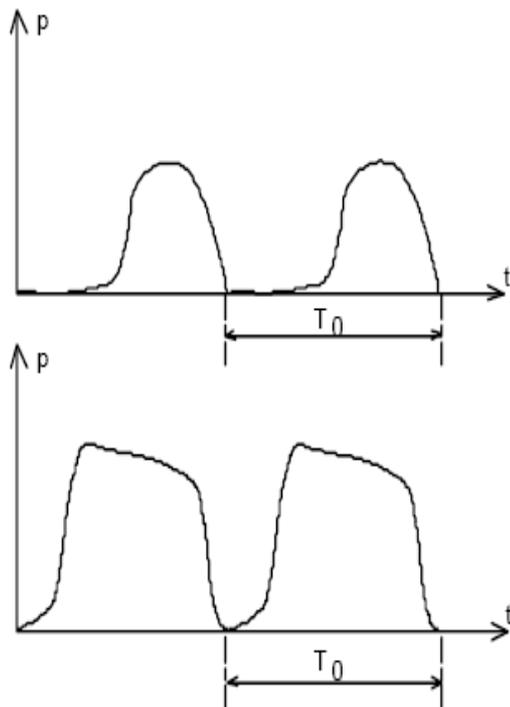
John Holmes' formant synthesizer (1964)



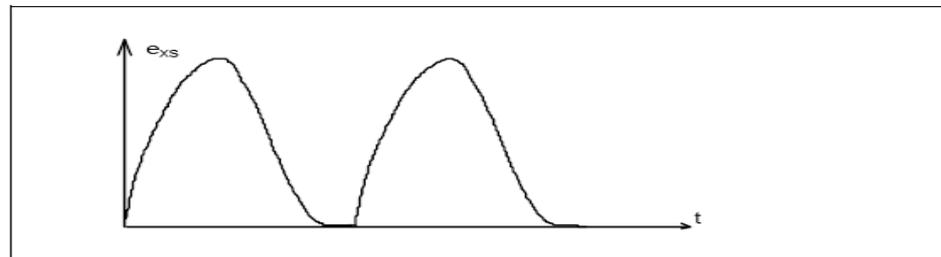
REALIZAREA EXCITATIEI

Reproducerea spectrului glotal se realizeaza cu diferite aproximații :

- aproximația triunghiulară
- aproximația polinomială
- excitarea unui filtru trece jos de ordinul doi, printr-un tren de impulsuri



1. Forme de undă glotice (Flanagan)



Undă glotică folosită de Fant

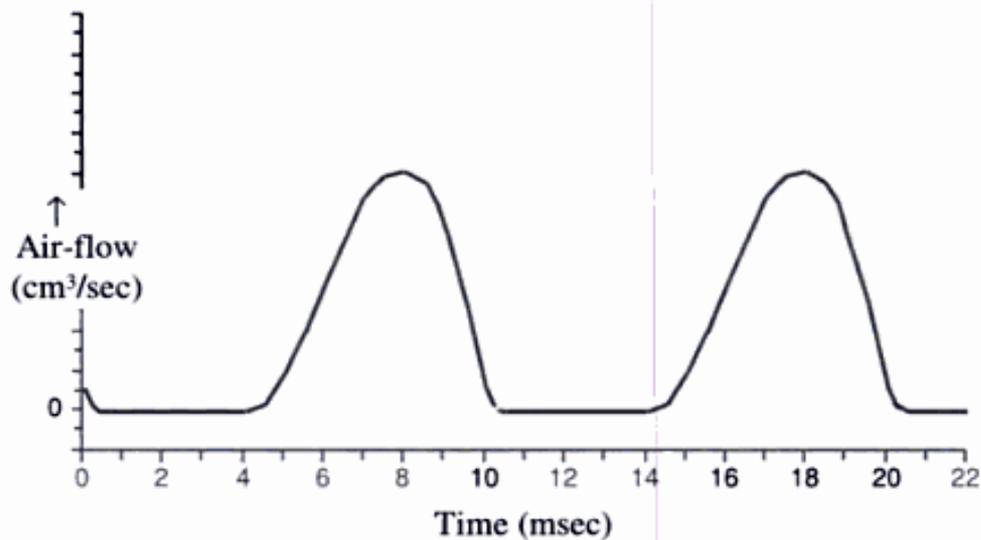


Figure 8.8 Two periods of an idealised glottal volume velocity wave with an F0 of 100 Hz

Speech acoustics

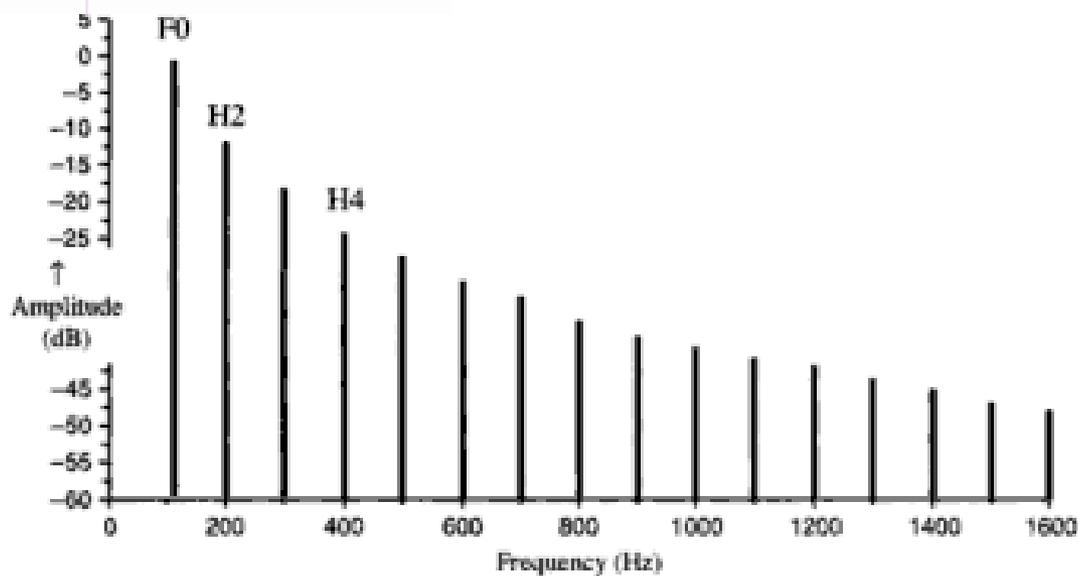
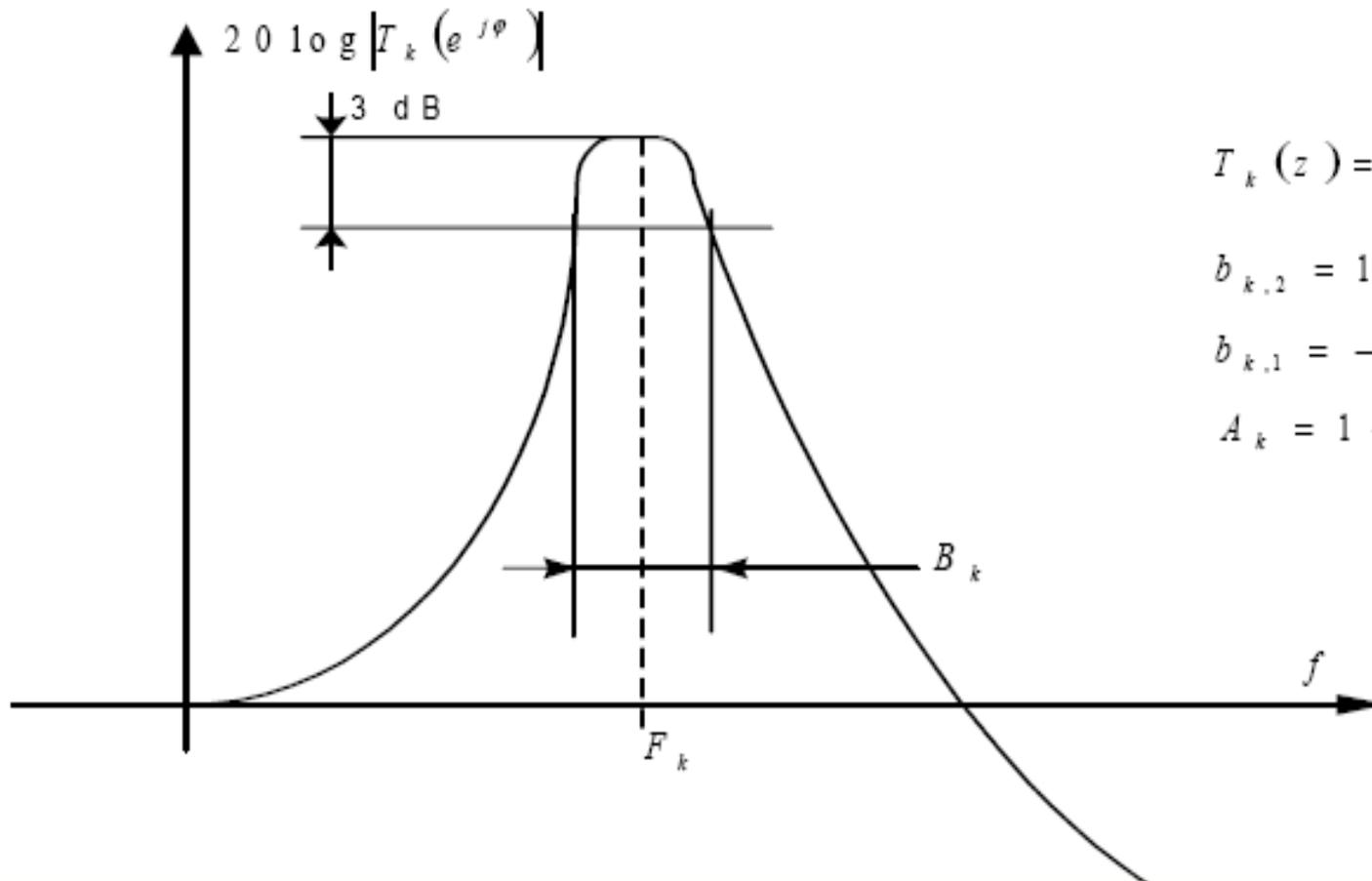


Figure 8.9 Idealised spectrum (up to 1.6 kHz) of the glottal volume velocity wave shown in Figure 8.8. The fundamental (F0) and second and fourth harmonics (H2, H4) are indicated

2.1. SINTEZA FORMANTICĂ - SINTEZA ÎN CASCADĂ

- Acest tip de sinteză constă în înserierea a 3-4 rezonatori de ordinul 2, care realizează fiecare, câte un formant de frecvență F_k și banda de trecere B_k
- Transmitanța se normalizează pentru a obține câștig nul la frecvență zero.



$$T_k(z) = \frac{A_k}{1 + b_{k,1}z^{-1} + b_{k,2}z^{-2}}$$

$$b_{k,2} = 1 - 2\pi B_k / f_c$$

$$b_{k,1} = -2b_{k,2} \cos 2\pi F_k / f_c$$

$$A_k = 1 + b_{k,1} + b_{k,2}$$

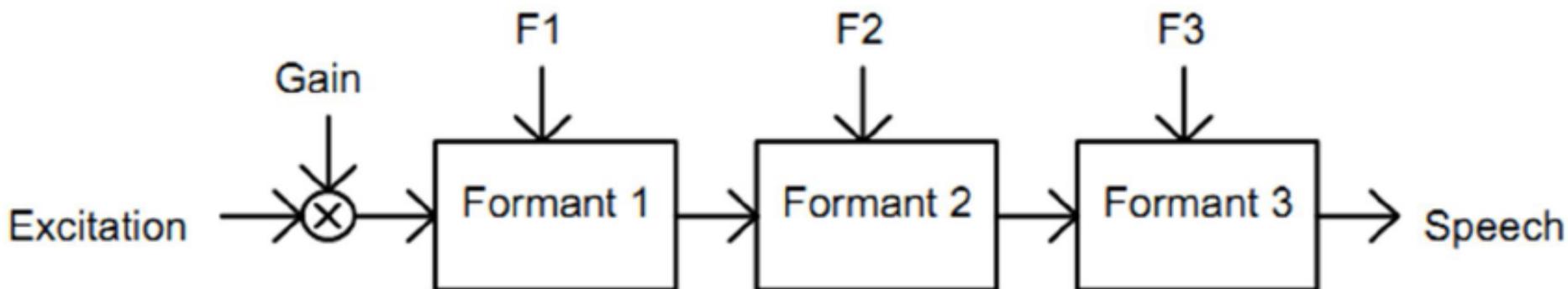
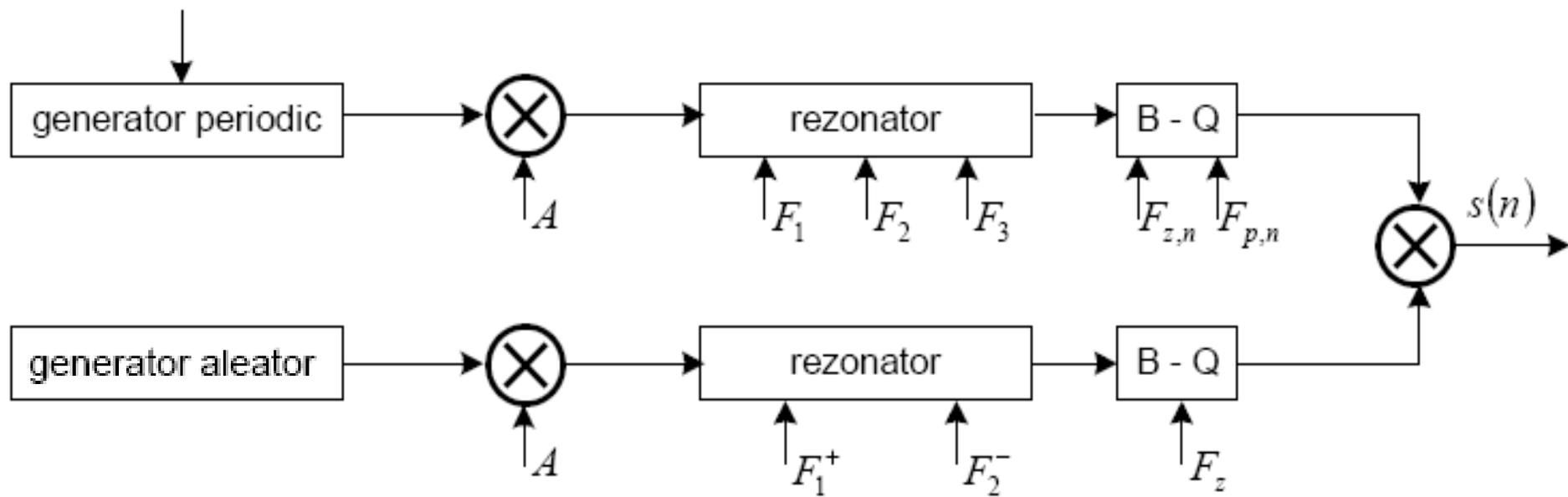


Fig. 3. Basic structure of cascade formant synthesizer [2]

Configurația serie - se bazează pe **funcția de transfer numai poli a tractului vocal**, ceea ce este **insuficient** pentru reprezentarea tuturor sunetelor vocii. Amplitudinea formanților superiori este funcție de a formanților inferiori, deoarece aceștia din urmă introduc în spectru o pantă de -12 dB/octavă. Acest fenomen nu se poate simula prin înserierea filtrelor.



$$V_k(z) = \frac{B_k}{1 + b_{1k} z^{-1} + b_{2k} z^{-2}}$$

$$b_{2k} = 1 - 2\pi B_k / f_{es}$$

$$b_{1k} = -2b_{2k} \cos(2\pi F_k / f_{es})$$

$$B_k = 1 + b_{1k} + b_{2k}$$

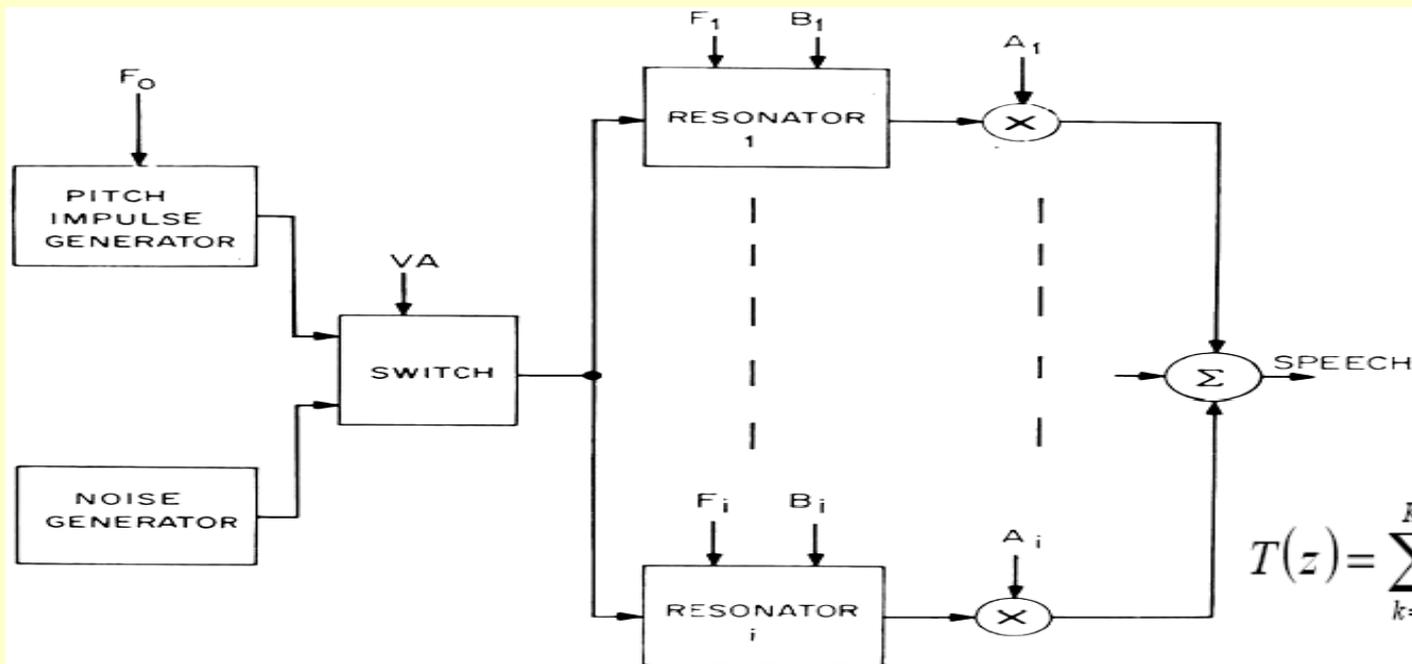
$$T(z) = \prod_{k=1}^K \frac{B_k}{1 + b_{1k} z^{-1} + b_{2k} z^{-2}}$$

	F (Hz)	B (Hz)
Sunete sonore		
Primul formant (F_1)	100 – 1100	45 – 130
Al doilea formant (F_2)	500 – 2500	5 – 190
Al treilea formant (F_3)	1500 – 3500	70 – 260
Rezonanță nazală ($F_{p,n}$)	200 – 1000	100
Antirezonanța ($F_{z,n}$)	200 – 1000	
Sunete nesonore		
Primul formant	200 – 500	60 – 300
Al doilea formant	1500 – 3500	60 – 200
Antiformant	800 – 2000	

- In tabel se indică domeniile de frecvențe ale parametrilor F_k și B_k (sunete sonore), $F_{z,n}$ și $F_{p,n}$ (sunete nazale) respectiv, și F_z (sunete fricative), valori determinate statistic

2.2. SINTEZA FORMANTICĂ - SINTEZA ÎN PARALEL

- Este mai delicată decât în cazul sintezei în cascadă, deoarece câștigul asociat fiecărui formant trebuie cunoscut
- Aceste câștiguri corespund rezidurilor polilor rezultați, din descompunerea transmitanței $T(z)$ în fracții simple
- Calcularea și realizarea lor trebuie să se facă cu mare precizie, fără a introduce zerouri în transmitanță



$$T(z) = \sum_{k=1}^K \frac{h_k}{1 + b_{1k} z^{-1} + b_{2k} z^{-2}}$$

Configurația paralelă - *filtre cu doi poli complex conjugați* este mai generală

- Se poate obține o funcție de transfer cu polii și zerourile dorite motiv pentru care această metodă este convenabilă pentru generarea sunetelor care conțin zerouri (de ex. sunetele nazale).

Avantaje

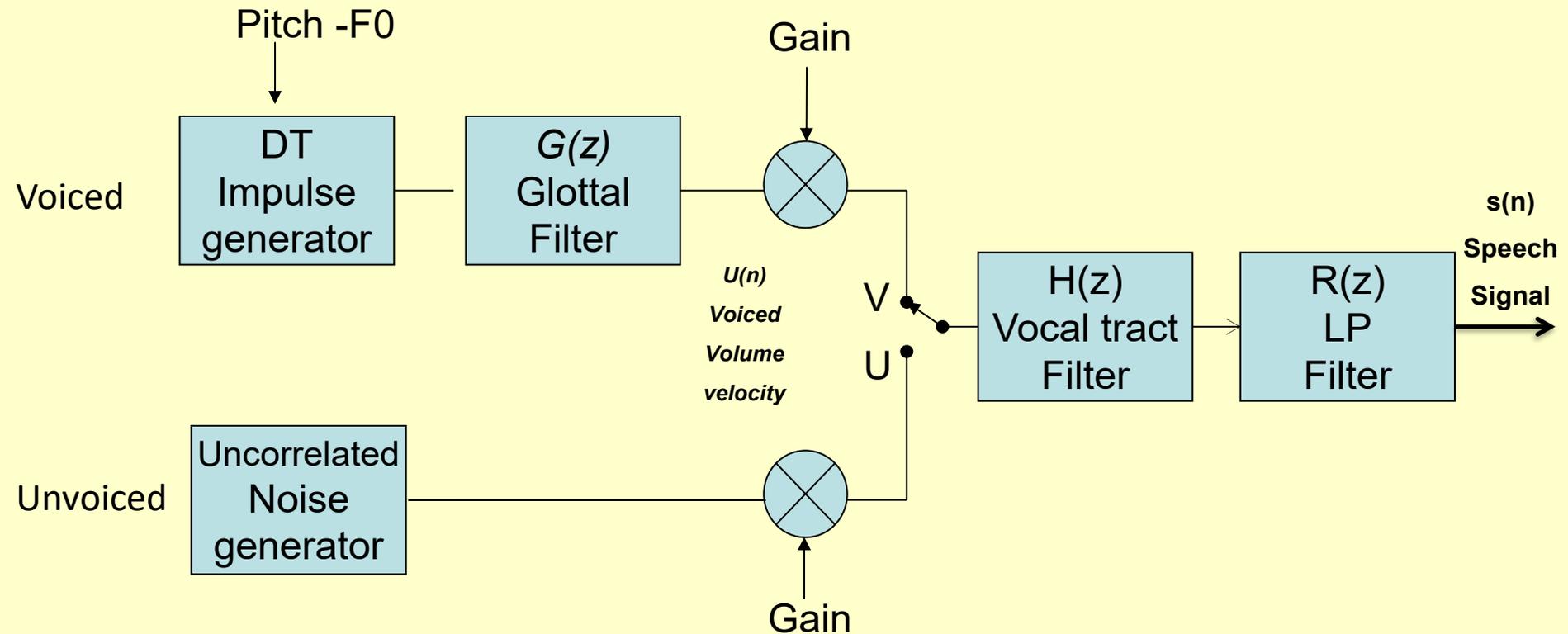
- Parametrii de sinteză sunt în *relație strânsă cu producerea și propagarea sunetului prin tractul vocal*
- Dacă se respectă condiția continuității în evoluția parametrilor, sintetizoarele formantice pot genera sunete sintetice cu *sonoritate plăcută*
- Pe lângă tipurile de excitație uzuale (sonor, nesonor) sintetizoarele formantice permit utilizarea *excitației mixte*

Dezavantaje

- Problema majoră a sintetizatoarelor formantice este **obținerea datelor** cu care se va opera
- *Urmărirea traiectoriei formanților (formant tracking)* este o sarcină foarte **dificilă** datorită faptului că această analiză nu poate fi automatizată complet
- Determinarea *lățimilor de bandă* aferente este și mai *problematică* (formanții sunt apropiați cu benzi suprapuse)

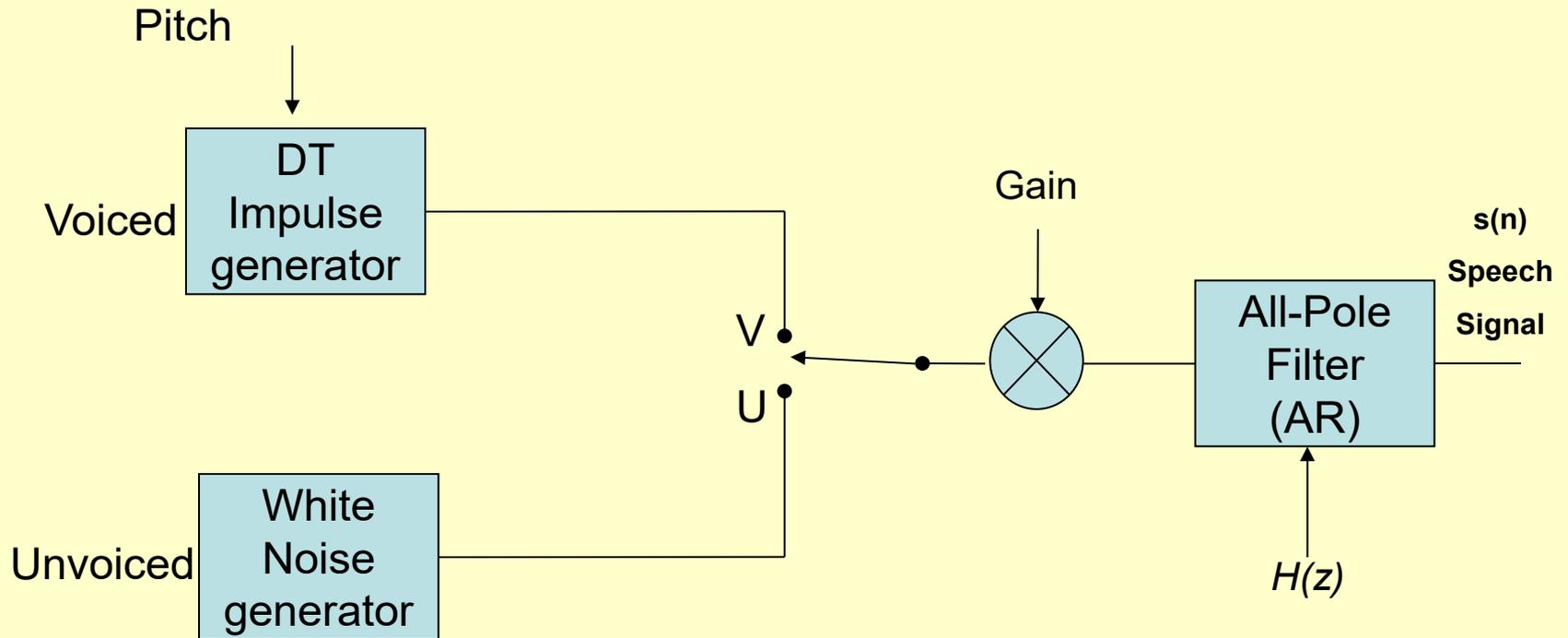
2.3. SINTEZA BAZATĂ PE PREDICȚIA LINIARĂ

- Modelul Real:



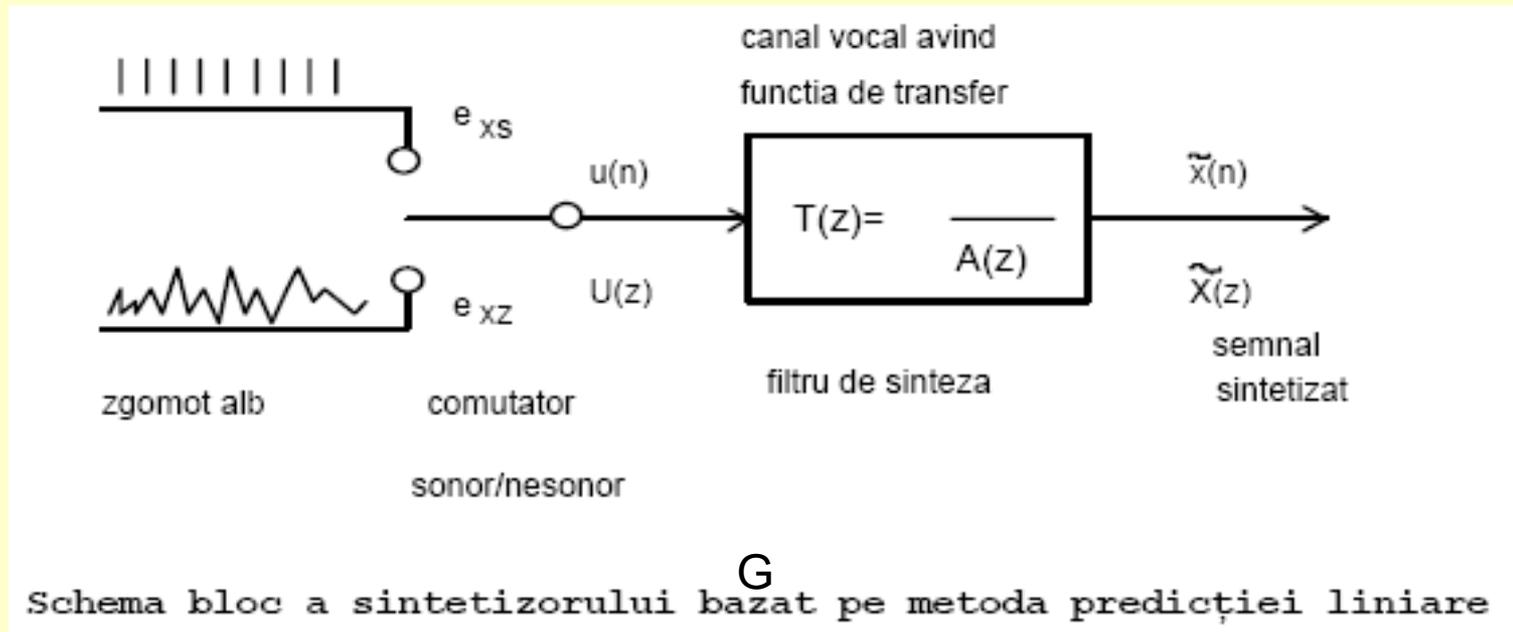
Modelul LPC real de producere a SV

- Modelul LPC simplificat utilizat in analiza



Modelul LPC

- are ca obiect estimarea parametrilor sistemului susceptibil de a genera un semnal artificial cât mai apropiat de semnalul original x .
- excitația se realizează fie cu un tren periodic de impulsuri, fie un zgomot alb
- pe lângă transmitanța canalului vocal se ia în considerare și forma reală de excitație (unda glotala)



Parametri modelului LPC:

- coeficienții filtrului numeric $A(z)$, p
- decizia sonor/nesonor
- F_0
- Câștigul, G

Decizia sonor/nesonor este un parametru necesar pentru care s-au elaborat diferite metode altele fata de cele uzuale:

1. metoda bazată pe transformata Fourier - Se știe că energia cadrelor sonore este concentrata în partea inferioară spectrului (0 - $F_{es}/4$), iar energia cadrelor nesonore este concentrat în partea superioara a spectrului ($F_{es}/4$ - $F_{es}/2$).

2. metoda bazată pe compararea energiei SV original cu a celui derivat. Experimental s-a constatat că *energia semnalului sonor original este mai mare decat a celui derivat*, iar în cazul semnalului nesonor acest raport este inversat.

3. metoda bazată pe NTZ a SV derivat. NTZ a semnalului derivat în cadrele nesonore a semnalului vocal este mai mare decât în cadrele sonore.

4. metoda bazată pe NTZ al derivatei funcției de autocorelație a SV original. NTZ al derivatei funcției de autocorelație a SV este considerabil mai mare în cazul semnalelor nesonore, decat în cazul semnalelor sonore.

Avantaje:

- Predicția liniară a devenit o metodă simplă de utilizat pentru simulare pe calculator
- Algoritmii utilizați pot fi ușor implementați hard
- Procesul de analiză este complet automatizat
- Calitatea vocii sintetice este destul de bună, în special din punctul de vedere al timbrului

Dezavantaje:

- decizia între excitația sonoră/nesonoră nu este întotdeauna satisfăcătoare
- Metodele noi (RELP, Modelul sursei mixte etc.) care și-au propus înlăturarea acestui neajuns, dar au tendința de a schimba timbrul vocii, deoarece excitația are un spectru plat mai întins

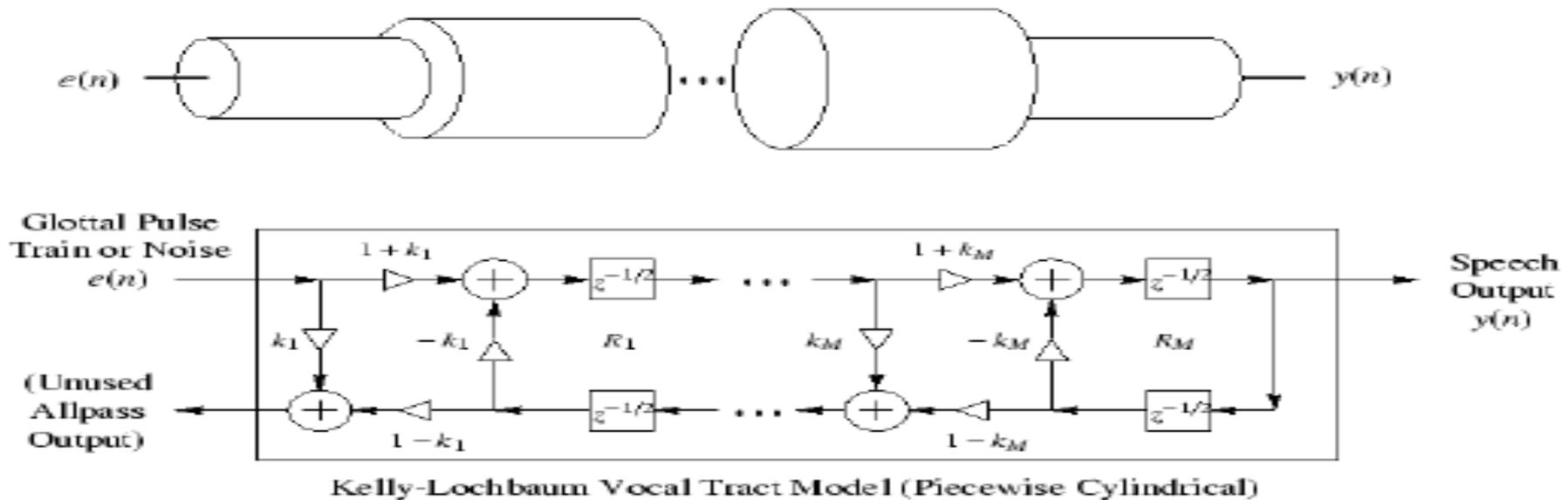
3. SISTEME DE SINTEZĂ A VOCII PRIN SIMULAREA TRACTULUI VOCAL

- **Principalul neajuns** al metodelor prezentate anterior este faptul că acestea *neglijează interacțiunea dintre sursă și tractul vocal*
- s-a constatat că diferite configurații ale tractului vocal pot produce sunete cu spectru identic
- altă abordare pentru descrierea acustică a tractului vocal este *simularea directă a generării și propagării undelor sonore* în interiorul sistemului fonator

- simularea acusticii tractului vocal poate fi combinată cu un model articulator
- sunt 2 metode capabile sa produca SV prin replicarea electro-acustica a mecanismului de producere a vorbirii:

- **METODA ANALOGIEI TRACTULUI VOCAL** - care simuleaza propagarea undei acustice prin tractul vocal

- **METODA ANALOGIEI TERMINALE** - care simuleaza structura spectrului de frecventa (caracteristicile rezonantelor si antirezonantelor care reproduc procesul de articulare)



3.1. METODA ANALOGIEI TRACTULUI VOCAL

- datele referitoare la forma tractului vocal sunt furnizate de *modelul articulato*r, respectiv de *modelul corzilor vocale*
- Modelele propuse pentru simularea vibrațiilor corzilor vocale:
 - **one-mass** de J. L. Flanagan și J. L. Landgraf;
 - **two-mass** de K. Ishizaka și J. L. Flanagan;
 - **multi-mass** de I. R. Tietze;
 - **two-beam** de R. Descot, J. Y. Aulage și B. Guerin
- *încearcă să calculeze cantitatea de aer glotal rezultată prin parametri ca: presiunea subglotală, tensiunea corziilor vocale și forma tractului vocal*
- *simularea propagării sunetului în interiorul tractului vocal în domeniul timp și în domeniul frecvență*

Punctul de plecare îl constituie **ecuațiile Webster**:

$$\frac{\delta^2 p}{\delta x^2} + \frac{1}{A} \frac{\delta p}{\delta x} \frac{\delta A}{\delta x} = \frac{1}{C^2} \frac{\delta^2 p}{\delta t^2} \qquad A \frac{\delta}{\delta x} \left(\frac{1}{A} \frac{\delta U}{\delta x} \right) = \frac{1}{C^2} \frac{\delta^2 U}{\delta t^2}$$

unde **p(x,t)** și **U(x,t)** sunt **presiunea**, respectiv **viteza volumică** în momentul **t** și poziția **x**, iar **A(x)** este aria secțiunii transversale a tractului vocal. ****

3.2. METODA ANALOGIEI TERMINALE

- simuleaza producerea SV printr-un circuit electric constand din cascada sau punerea in paralel a unor circuite care descriu rezonantele (formantii) si antirezonantele (antiformantii)
- F_k si B_k ale fiecarui circuit sunt variabile (metoda Formantica)
- conexiunea in cascada este avantajoasa in stabilirea automata a raportului intre amplitudinile formant/antiformant (vocale-structura spectrala)
- conexiunea paralela e avantajoasa ptr ca forma spectrala finala poate fi mai precis simulata (nazale, fricative)

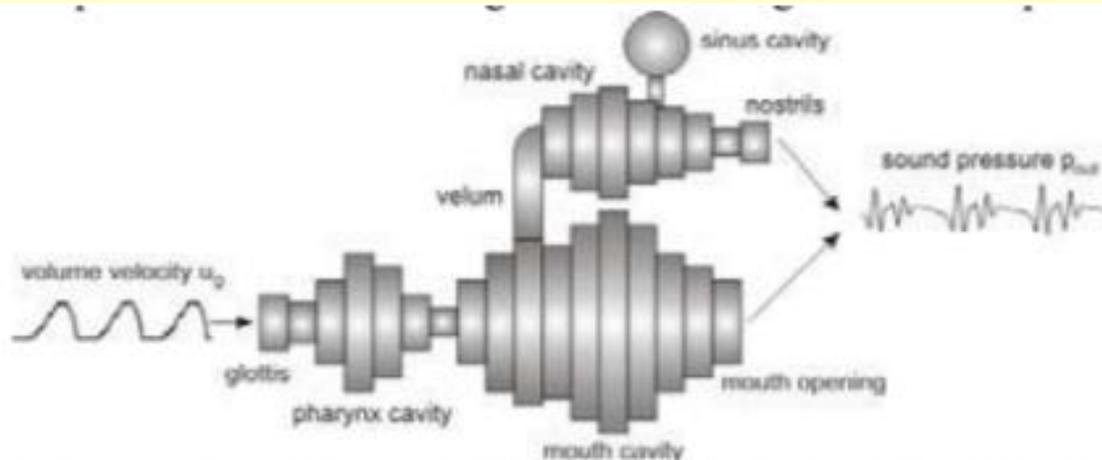


Fig. 2: An articulatory model of the human speech production system

4. SINTEZA PE BAZĂ DE REGULI

- se construiește vocea prin reguli fonetice, fonologice și fonetico-prosodice derivând direct din descrieri lingvistice ale limbii țintă. (pe baza unor secvențe de simboluri *fonetice / silabice / text*)
- descrierea fonetică a textului (foneme, reguli fonetice), reguli de pronunție, reguli de transformare din text în formă canonică (etalon)
- modele de prosodie: reguli pentru intonație, creșterea/reducerea intensității (accentuare), durata fonemelor, pauzele.
- interpretabilitate ridicată: parametrii sunt în mod explicit controlați, facilitând studierea relațiilor între text și vorbire.
- adecvate pentru vocabular limitat sau când este nevoie de control strict asupra caracteristicilor vocale.
- naturalitatea și eficiența la vocabulare mari sunt adesea limitate; adaptarea la dialecte/limbi diferite necesită efort semnificativ în definirea regulilor.
- caracteristicile *unităților fundamentale* de vorbire (silabe, foneme sau vorbirea pe T0) sunt stocate și se concatenează pe baza unor reguli
- Caracteristicile prosodice (F0, A,..) sunt de asemenea controlate prin reguli
- *Calitatea unităților fundamentale* pentru sinteză ca și *regulile de control* joacă un rol esențial în această metodă și trebuie să se bazeze pe caracteristicile fonetice și lingvistice ale vorbirii naturale specific limbii respective
- În cazul sistemelor *bazate pe forma de undă*: dacă unitățile fundamentale pentru sinteză sunt fonemele (extrase din vorbirea naturală sau artificială) memoria necesară stocării este foarte mică (de obicei sunt 30-50 de foneme), dar regulile de conectare ale fonemelor sunt foarte complicate, iar calitatea sintezei este redusă
- se utilizează frecvent unități mai mari decât *fonii sau alofonii* (foneme dependente de context) - pentru o sinteză de calitate sunt necesare alte unități fundamentale $\sim 10^3-10^4$

Table 16.4 Unit types in English assuming a phone set of 42 phonemes. Longer units produce higher quality at the expense of more storage. The number of units is generally below the absolute maximum in theory: i.e., out of the $42^3 = 74,088$ possible triphones, only about 30,000 occur in practice.

Unit length	Unit type	#Units	Quality
Short  Long	Phoneme	42	Low  High
	Diphone	~1500	
	Triphone	~30K	
	Demisyllable	~2000	
	Syllable	~11K	
	Word	100K–1.5M	
	Phrase	∞	
	Sentence	∞	

Alegerea unitatilor

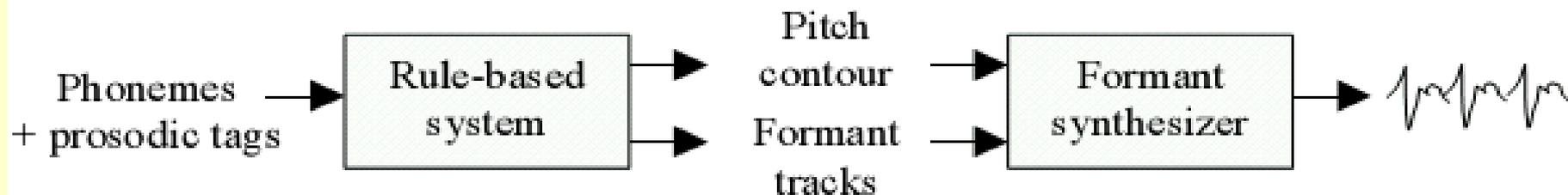


Figure 16.2 Block diagram of a synthesis-by-rule system. Pitch and formants are listed as the only parameters of the synthesizer for convenience. In practice, such system has about 40 parameters.

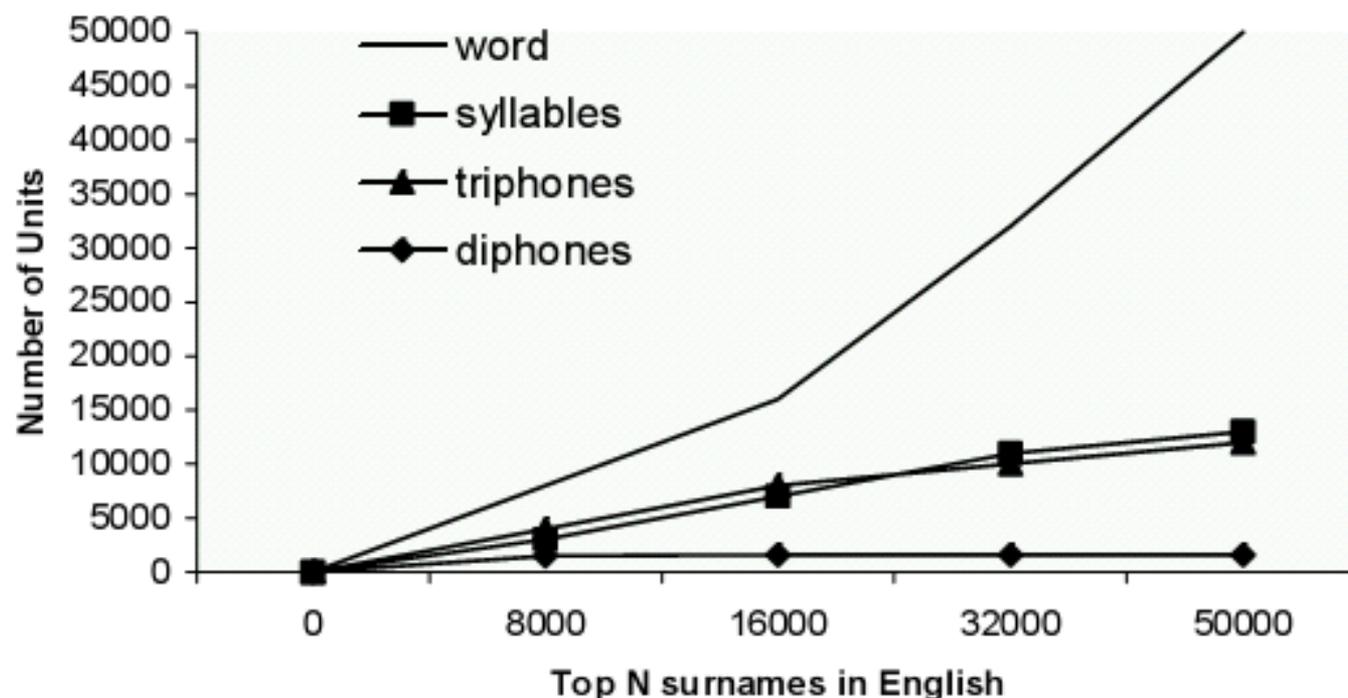


Figure 16.5 Coverage with different number of units displays the number of units of different types required to generate the top N surnames in the United States [34].

CONCATENAREA UNITATILOR fundamentale

- CUVINTE

- Nu are acoperire completă pentru domenii largi

- Capacitatea limitată de a modifica F0, amplitudinea și durata fără a pierde din naturalețe și inteligibilitatea

- Necesitatea unei baze de date imense pentru a extrage mai multe versiuni ale fiecărui cuvânt

- SUB-UNITATI (silabe, demisilabe, difoni,...)

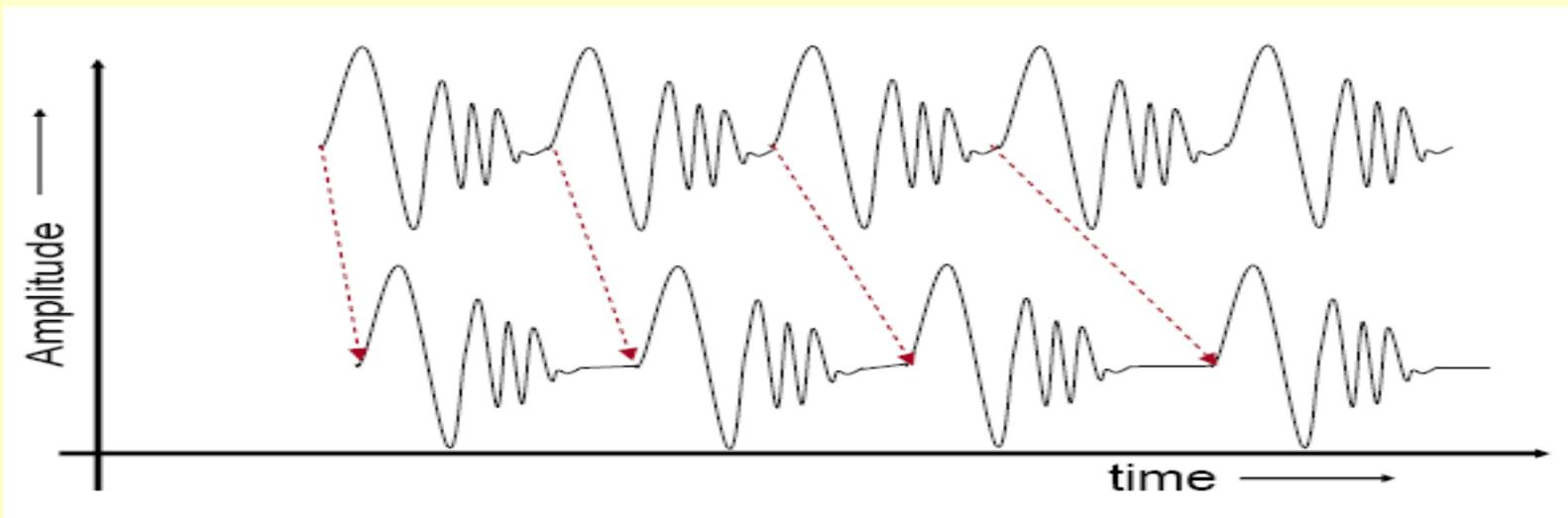
- *greu de izolat în context din cauza coarticularii*

- necesita variante alofonice pentru a caracteriza subunitățile în toate contextele in care apar

- se cere efort mare de prelucrare a semnalului pentru a netezi conexiunile cu unitatile adiacente

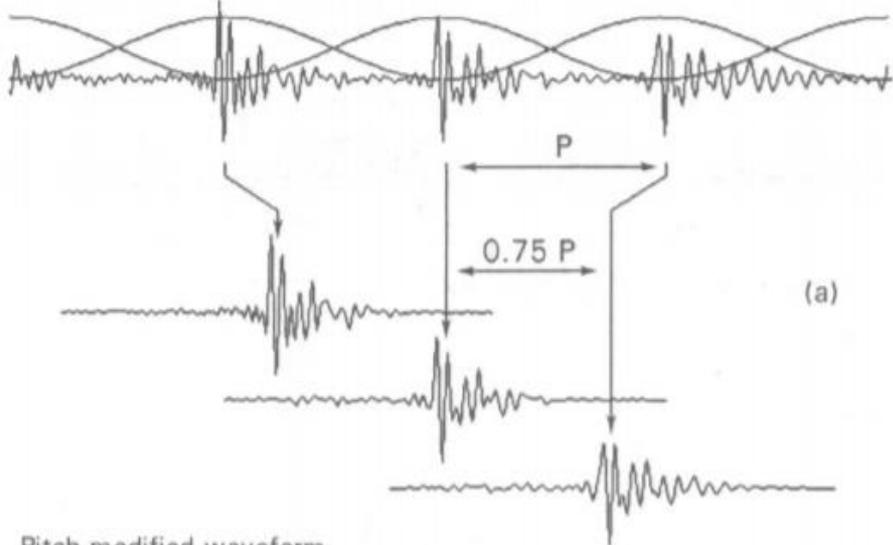
CONCATENAREA UNITATILOR REPREZENTATE

- **LPC**
 - Simplu, ușor de concatenat unitățile, eficient pentru modificarea F0
 - Nu funcționează bine pentru sunete nazale (lipsa zerouri nazale)
 - Excitație glotală nu este corectă (excitația asumată ca impulsuri)
 - LPC - nu funcționează bine pentru excitații mixte
- **TD-PSOLA—Time Domain, Pitch Synchronous Overlap Add Synthesis** -
Modificarea eficientă a prozodiei (sincron cu F0)
 - Nu cere o uniformizare la punctele de jonctiune



Pitch-scaling

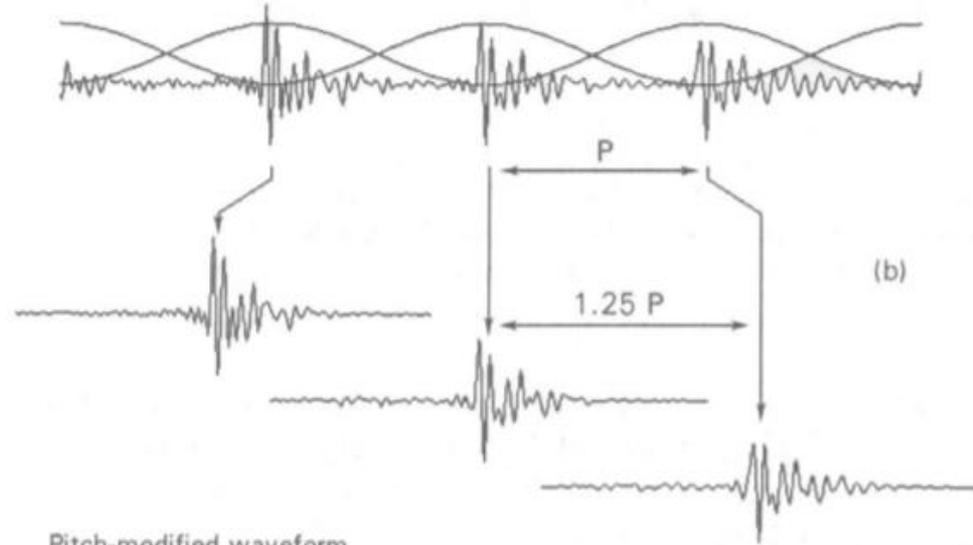
Original waveform



Pitch-modified waveform



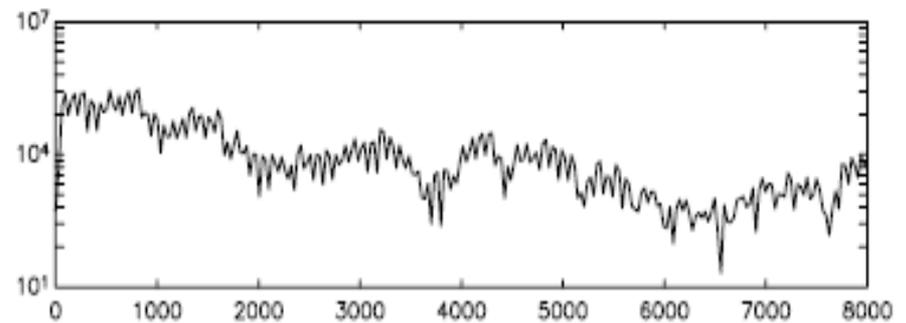
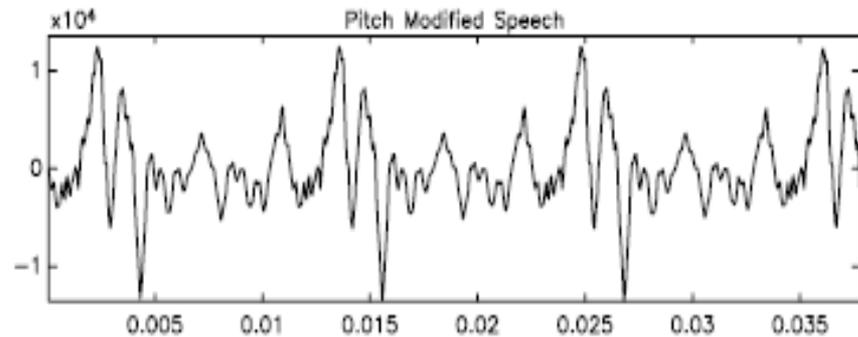
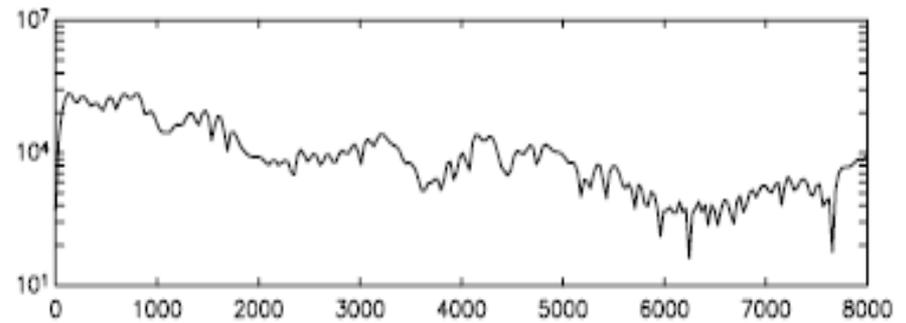
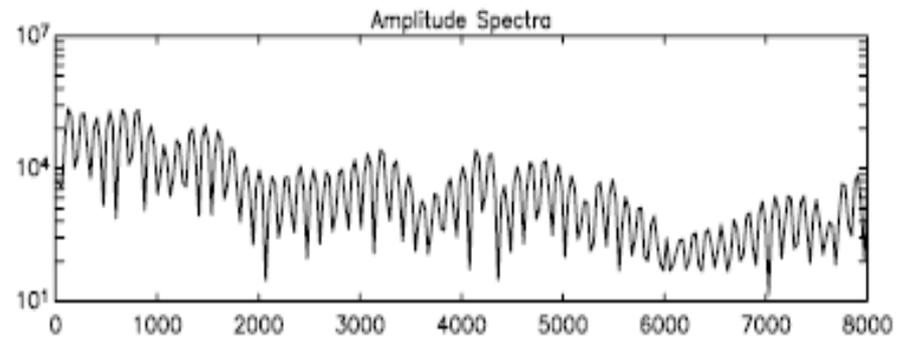
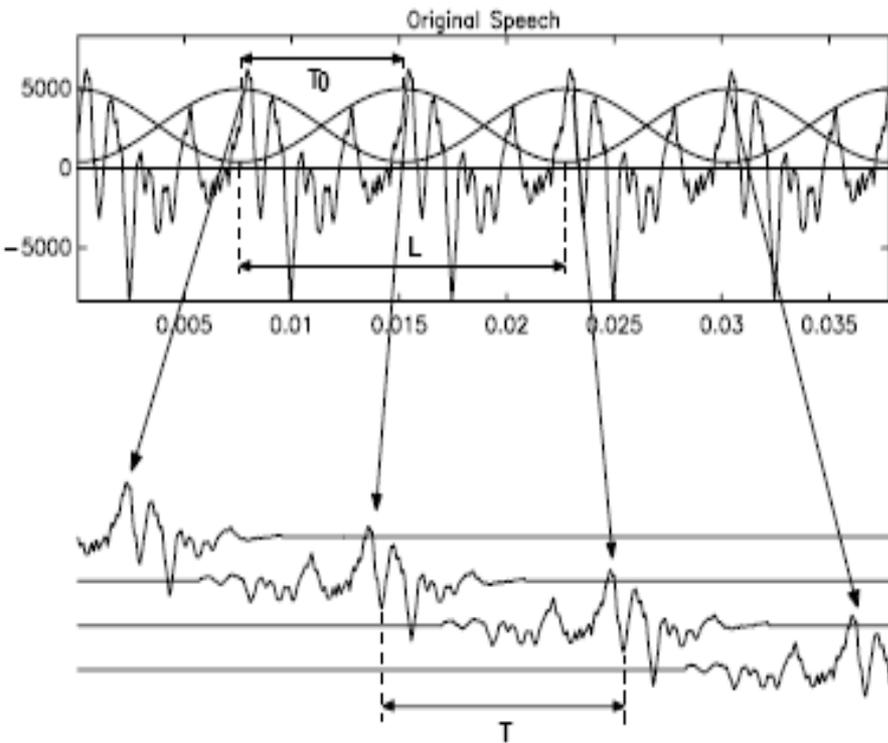
Original waveform



Pitch-modified waveform



SINTEZA TD-PSOLA



(MBROLA- The *Multi-Band Re-synthesis OverLap-Add*)

- sinteza in domeniul timp care combina avantajele PSOLA (calcul redus, modificarea prozodieii), cu un algoritm hibrid timp-frecvență, off-line pentru uniformizarea tranzițiilor.

Exemplu.

Pentru **limba japoneză** este nevoie de ~200 de silabe de tip CV/VC (consoană–vocală) sau 5000-6000 de unități de tip CVC (putând fi reduse la ~1000) sau între 700-800 de unități de tip VCV.

De exemplu cuvântul SAKURA = cireș poate fi reprezentat cu unități :

CV : SA+KU+RA

CVC: SAK+KUR+RA

VCV: SA+AKU+URA.

Unitățile CVC sunt conectate pe consoană, iar cele VCV pe vocală fiecare prezentând anumite avantaje în ușurința de conectare.

Pentru **limba engleză** sunt > de 3500 de silabe care se extind la 10.000 dacă se consideră și alofonii.

-silabele se descompun în diade/difoni (400-1000) sau demisilabe (~1000) care au reguli mai simple de concatenare.



Length	Unit	# Units (English)	# Rules, Necessary Unit Modifications	
Short			Many	Quality Low
	Allophone	60-80		
	Diphone	$< 40^2 - 65^2$		
	Triphone	$< 40^3 - 65^3$		
	Demisyllable	2K		
	Syllable	11K		
	VC*V			
	2-syllable	$< 11 \text{ K}^2$		
	Word	100K-1.5M		
	Phrase	∞		
	Long	Sentence		

5. Metode neuronale end-to-end actuale

Iată câteva din cercetările și abordările populare și actuale ale sintezei vocale:

- [WaveNet: A Generative Model for Raw Audio](#)
- [Tacotron: Towards End-to-End Speech Synthesis](#)
- [Deep Voice 1: Real-time Neural Text-to-Speech](#)
- [Deep Voice 2: Multi-Speaker Neural Text-to-Speech](#)
- [Deep Voice 3: Scaling Text-to-speech With Convolutional Sequence Learning](#)
- [Parallel WaveNet: Fast High-Fidelity Speech Synthesis](#)
- [Neural Voice Cloning with a Few Samples](#)
- [VoiceLoop: Voice Fitting and Synthesis via A Phonological Loop](#)
- [Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions](#)

Tema-studiu individual.

- Noile abordari ale sintezei SV:
 - (**Wavenet**) utilizarea rețelelor neuronale, modelul folosit este probabilistic și autoregresiv
 - (**Tacotron**) este un model generativ text-vorbire, care sintetizează vorbirea direct din perechi de text și audio. Tacotron generează vorbire la nivel de cadru și, prin urmare, este mai rapid decât metodele autoregresive la nivel de eșantion.
 - (**Deep Voice 1-3**) este un sistem text-to-speech dezvoltat folosind deep neural networks, DNN.
 - (**Parallel WaveNet**) Fast High-Fidelity Speech Synthesis - o metodă cunoscută sub denumirea de Probability Density Distillation, care antrenează o rețea paralelă feed-forward alimentată de la un WaveNet antrenat.

<https://www.deepmind.com/blog/>

<https://heartbeat.fritz.ai/a-2019-guide-to-speech-synthesis-with-deep-learning-630afcafb9dd>

Stuidu.

<https://www.techradar.com/best/best-text-to-speech-software>

<https://heartbeat.fritz.ai/a-2019-guide-to-speech-synthesis-with-deep-learning-630afcafb9dd>

