

Sesión 1. Psicoacústica de la percepción espacial del sonido

• Introducción

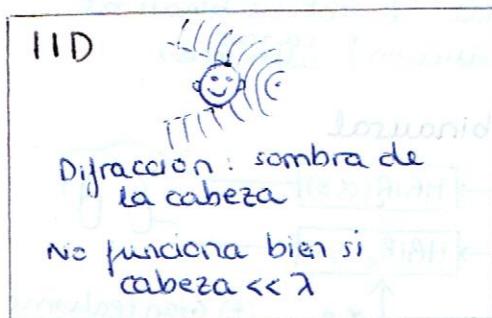
- El modelo cerebral del oído es más complejo que el visual
→ localiza con rapidez y precisión el origen del sonido

• Sensibilidad del oído



• Localización espacial

• Diferencia de intensidad



• Localización espacial en el plano horizontal

• Diferencia de tiempo interaural (ITD)

ITD



$$\text{ITD} = \frac{c}{v} (\theta + \text{seno})$$

La ITD se interpreta como una diferencia de fases
(el oído detecta fase, no retardo)

Ambigüedad de fase a altas frecuencias

confusión delante-detrás

Fuente 1 misma ITD

Fuente 2 cono de confusión



$\rightarrow f > 1.5 \text{ kHz} \Rightarrow$ Ya no sirve ITD

Mecanismo combinado:

↳ transición gradual

$$\left\{ \begin{array}{l} \theta < 1.5 \text{ kHz} \rightarrow \text{ITD} \text{ (detecta desfase)} \\ \theta > 1.5 \text{ kHz} \rightarrow \text{IID} \text{ (difracción cabeza)} \end{array} \right.$$

• Localización espacial en elevación

- Filtrado de las orejas y hombro según la dirección
↳ ecos con retardos

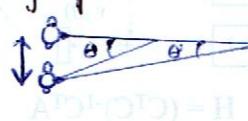
Ambos oídos están en el mismo plano (no se puede detectar desfase en elevación)

• Percepción de la distancia

→ Intensidad: $A_t = \frac{P_t}{P_r}$ física: propagación atenua la distancia
inteligencia: según el tipo de sonido (voz humana, pájaro, ...) el cerebro estima P_t

→ Atenuación de altas frecuencias: el aire tiene mas 'rozamiento' a altas frecuencias

→ Paralaje por movimiento: - según la relación entre faltas y fijas el cerebro deduce la distancia



se mueve ligeramente la cabeza y se compara el ángulo (típico para encontrar a un grialo)
(así es como se mide la distancia de las estrellas)

→ Exceso de diferencia de IID: para objetos muy cercanos → ej: mosquito en el oído

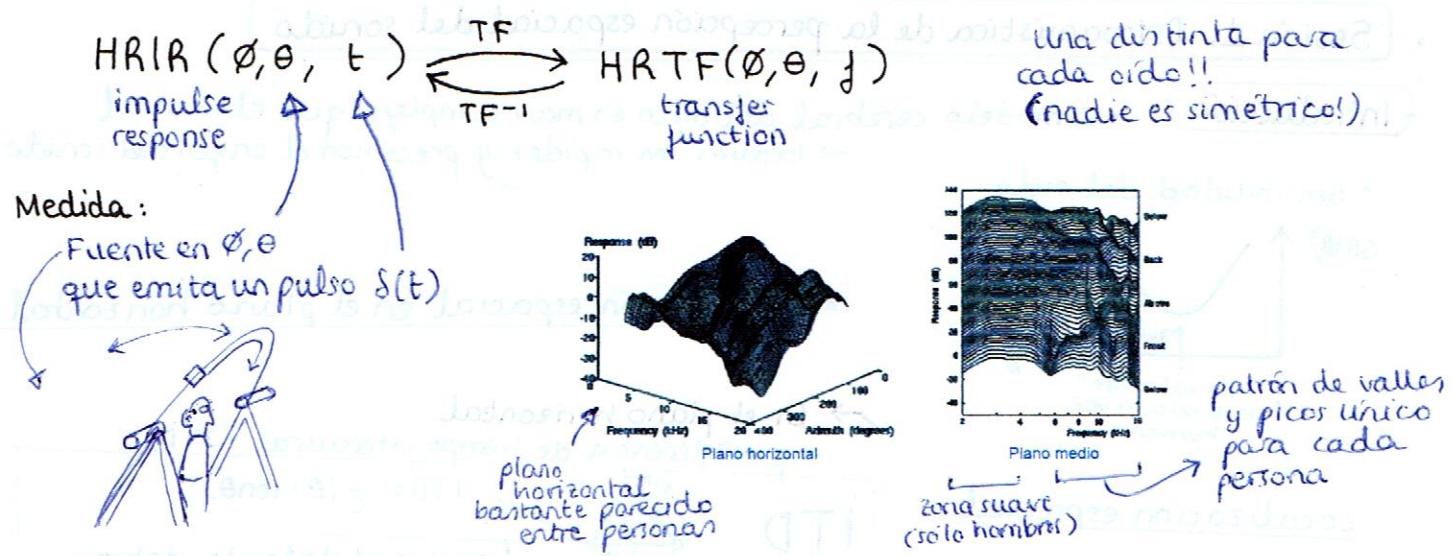
→ Relación entre sonido directo y reverberación: (en salas)



la energía del sonido directo depende de la distancia

la energía TOTAL de los ecos es más o menos igual en todos los puntos de la sala

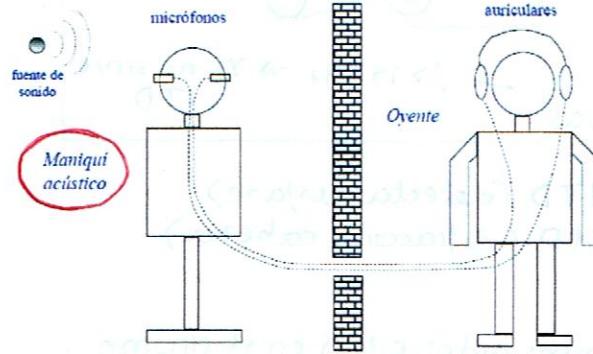
Función de transferencia relacionada con la cabeza (HRTF)



Reproducción del sonido espacial

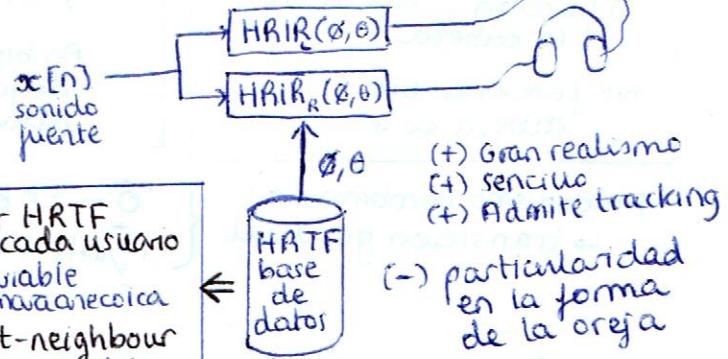
Sistemas basados en la HRTF

Grabación/reproducción binaural



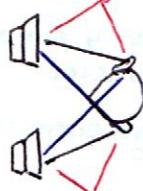
obtener la { - grabación binaural
señal - síntesis binaural
A reproducción {- carcelas
altavoces}

Síntesis binaural



Reproducción binaural con altavoces:

Problema
- crosstalk
- reflexiones

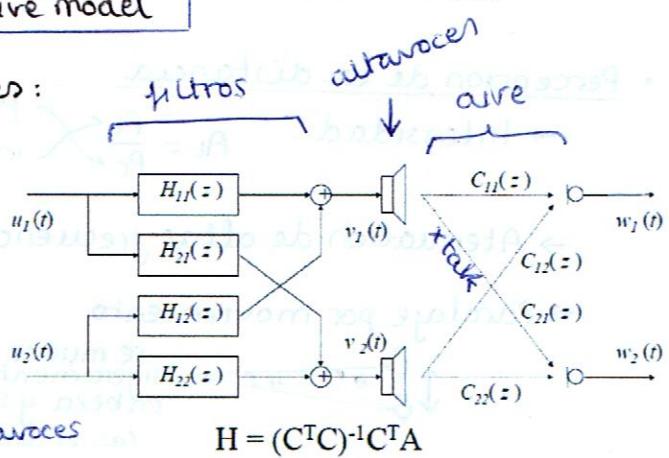


Solución: cancelar los caninos cruzados

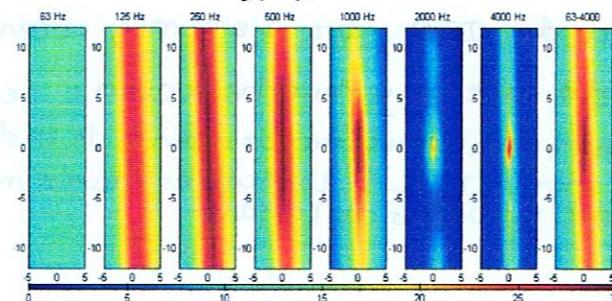
↳ sólo sirve si el usuario no se mueve del sweet spot: estudios

- teóricos en espacio libre
- medidas

altavoces



movearse delante/detrás
no afecta mucho



Degradación (relación directa / interferente) en función del desplazamiento

En 2 kHz hay singularidad debido a distancia entre oídos → único sweet spot

EVOLUCIÓN DEL SONIDO ENVOLVENTE

• Estéreo



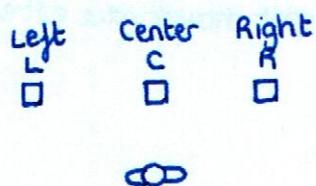
Efecto phantom:

- Por diferencia de volumen el cerebro interpreta una posición entre L y R
- Problema: no hay diferencia entre tiempos de llegada
- Solución: el cerebro se acostumbra
- Problema 2: sólo funciona si estás centrado

potenciómetro PAN-POT (panorama potentiometer)
mci: L R

$$\begin{aligned} L &= y \cdot V \\ R &= x \cdot V \end{aligned}$$

• Canal central para diálogo

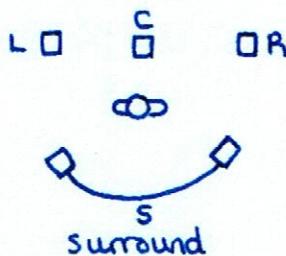


Los espectadores no centrados se desorientaban al oír el diálogo a los lados.

Canal central para el diálogo

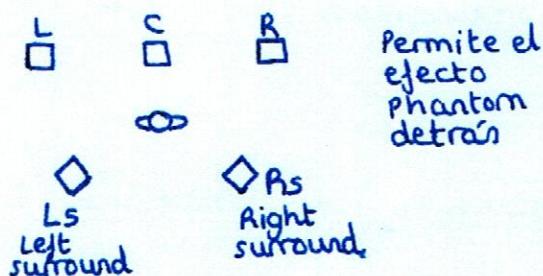


• Canal surround



- Llena el ambiente
- Simula reverberación
- Efectos/surtos

• Desdoblar canal surround



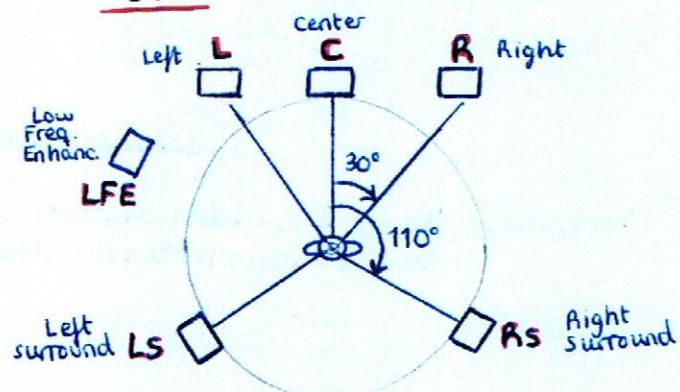
Permite el efecto phantom detrás

• LFE: low frequency enhancement

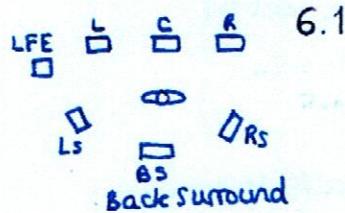
usa gran altavoz (subwoofer) para frecuencias [20Hz, 120Hz]

No requiere posicionamiento ya que apenas hay diferencia de fase y volumen entre los oídos a bajas frecuencias (\rightarrow difracción de cabeza \rightarrow)

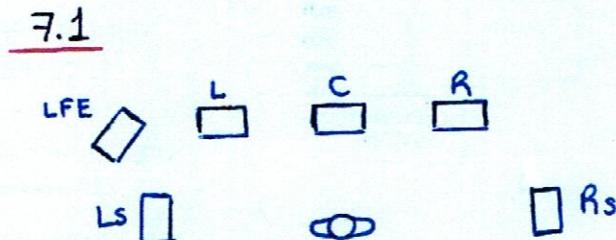
5.1 (estándar ITU)



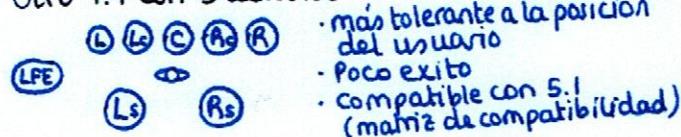
• Tercer canal de surround



• Desdoblar el Back surround \rightarrow 7.1



Otro 7.1 con 5 delante



Stereo: - phantom image created

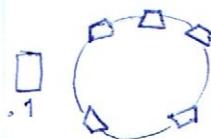
X-Y technique: variando intensidad ← Pan pot: el más usado
A-B technique: variando time-of-arrival

Stereo Sweet Spot:



Surround 5.1

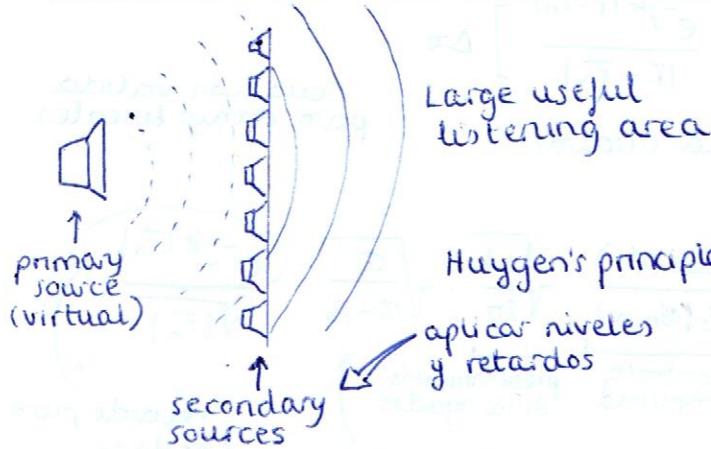
- pan-pot delante con 3 altavoces → sweet spot mayor que estéreo
 - localización lateral y debajo es pobre



Solución → Wave Field Synthesis

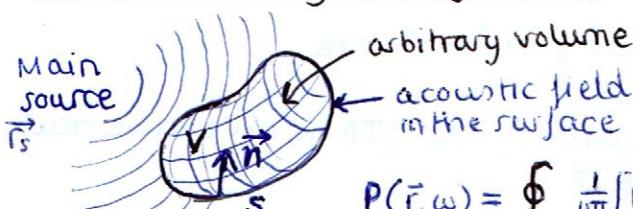
SESIÓN 2. WAVE-FIELD SYNTHESIS

→ 2'5D (no tiene elevación)



Otra técnica aún mejor:
Ambisonics: descomposición del campo acústico en armónicos esféricos (con array de 32 micrófonos se capta hasta 4º orden)

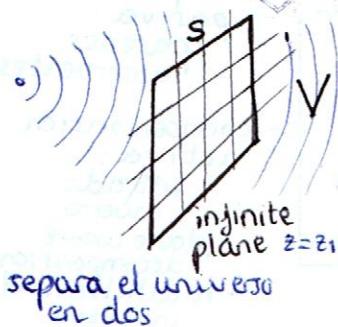
WFS : Theory Background



$$P(\vec{r}, \omega) = \oint_S \frac{1}{4\pi} \left[P(\vec{r}_s, \omega) \cdot \frac{\partial}{\partial n} \left(\frac{e^{-jk|\vec{r}-\vec{r}_s|}}{|\vec{r}-\vec{r}_s|} \right) + \frac{\partial P(\vec{r}_s, \omega)}{\partial n} \left(\frac{e^{-jk|\vec{r}-\vec{r}_s|}}{|\vec{r}-\vec{r}_s|} \right) \right] dS$$

presiones = monopolar
= altavoces

derivada de presiones = dipolos



Rayleigh integral equation

$$P(\vec{r}, \omega) = |z - z_1| \int_S \left[P(\vec{r}_s, \omega) \frac{1 + jk|\vec{r} - \vec{r}_s|}{2\pi|\vec{r} - \vec{r}_s|^3} e^{-jkl|\vec{r} - \vec{r}_s|} \right] dS$$

↓ Discretizar las fuentes

↑ Ya no necesito dipolos!

$$P(\vec{r}, \omega) = |z - z_1| \sum_n \left[P(\vec{r}_n, \omega) \frac{1 + jk|\vec{r} - \vec{r}_n|}{2\pi|\vec{r} - \vec{r}_n|^3} e^{-jkl|\vec{r} - \vec{r}_n|} \right] \Delta x \Delta y$$

↓ plano infinito de altavoces

Simplification to a line array

↓ suponer que cada altavoz emite onda cilíndrica en lugar de esférica (no es cierto, pero el error es pequeño si estamos en el plano de los altavoces)

$$P(\vec{r}, \omega) = jkpc \sum_n u_n(\vec{r}_n, \omega) \frac{e^{-jkl|\vec{r} - \vec{r}_n|}}{|\vec{r} - \vec{r}_n|} \Delta x$$



(aparece como rotación horizontal 180°)
↓ si las ondas están en la misma dirección
↓ en la mitad de distancia

(rotación sobre el eje vertical) JMF 2005
↓ en la mitad de distancia en la otra dirección

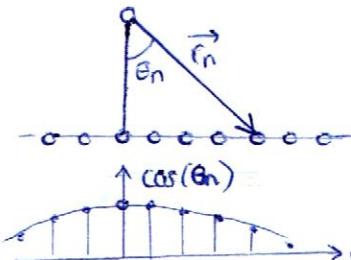
↓ se cancelan las ondas con el efecto de la difusión de la onda en el espacio -
se reduce la amplitud en un 1/4 (JMF 2005)

JMF 2005

Para N altavoces

$$P(\vec{r}, \omega) = \sum_{n=1}^N \left[Q(\vec{r}_n, \omega) \cdot G(\phi_n, \omega) \cdot \frac{e^{-jk|\vec{r} - \vec{r}_n|}}{|\vec{r} - \vec{r}_n|} \right] \Delta x$$

Excitación de cada altavoz



$$Q(\vec{r}_n, \omega) = S(\omega) \cdot \frac{\cos(\phi_n)}{G(\phi_n, \omega)}$$

$$\frac{\sqrt{jk}}{2\pi} \cdot \sqrt{\frac{\sigma_0}{\sigma_0 + \rho_0}}$$

$$\frac{e^{-jk|\vec{r}_n|}}{\sqrt{|\vec{r}_n|}}$$

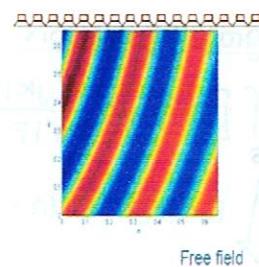
retardo puro
(dejarse variable con la frecuencia)

Ecuación válida para arrays lineales

Limitaciones:

- Muchos altavoces: cada uno con su señal y su hardware
- Limitación: no se puede jugar con elevación?
↳ ¿Y si hacemos phantom effect en elevación? → No, IID no existe en vertical
↳ mejor solución: combinar con HRTF
- Tamaño de array limitado (truncation)
- Efecto de la sala

Ecos estropean el efecto

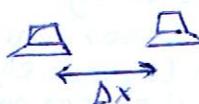


Soluciones:

- compensación pasiva:
paredes absorbentes
- compensación activa:
filtrado inverso
- plane wave decomposition
- multichannel inversion

→ Spatial aliasing

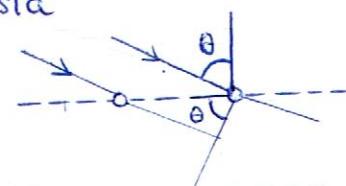
. separación entre altavoces



. sólo se puede sintetizar bien hasta

$$\Delta x \cos \theta_{\max} < \frac{\lambda_a}{2}$$

$$f_a = \frac{c}{2\Delta x \cos \theta_{\max}}$$



Soluciones:

- ignorar el problema → en el ITEAM $\Delta x = 18 \text{ cm}$, $f_a = 1 \text{ kHz}$
la fase errónea solo se nota si corren lateralmente y mirando hacia el array

· usar más altavoces agudos

· OPSI (Optimized Phantom Source Image)

↳ Usar WFS para $f < f_a$

↳ Usar efecto phantom para $f > f_a$

· usar DML (Distributed Mode Loudspeaker)

↳ membrana que se mueve entera con MAP (multi actuator)

- array de 32 (4x8)

- array de 96 (12x8) → de los más grandes de Europa

- MAP(DML) → 4 paneles x 6 excitadores

WFS in GTAC

Applications of WFS

- Telepresencia
- Realidad virtual
- Reproducción de música con fidelidad espacial
- Planetarium, IMAX, ... (short term)
movie theaters (medium term)
home cinema (long term)

Research lines for WFS

- | | |
|--|---|
| <ul style="list-style-type: none"> - room compensation - Discrete Time Modelling (FDTD) - Authoring & Real Time Tools | <ul style="list-style-type: none"> - High-power applications - improvements in DMLs |
|--|---|

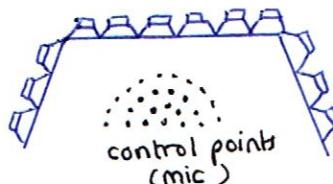
Compensación de la sala

- Plane-wave decomposition →



emitir los ecos "al revés"

- MIMO inverse filters



- measure impulse response between each loudspeaker and each control point.
- obtain bank of inverse filters
- Simulate sound field

corrección multipunto → se han obtenido buenas correcciones en un gran área

Simulación FDTD

- Microfonos virtuales para obtener RIR
- Superordenadores

High-power applications

- Diameter of loudspeakers
- Directivity

Experiments & simulations :

- Field analysis
- Precise aliasing Frequencies

Authoring tools for WFS

- standalone application
- VST plugin that works with 'Cubase'

Software Block Diagram

SESIÓN 3. SOUND SOURCE SEPARATION (SSS)

III - 1

INTRODUCTION

- Recovering each source signal from a given mixture
- Difícil → active research

Applications :

- remixing
- speech enhancement (e.g. hearing devices)
- speech recognition
 - ↳ singer recognition and melody extraction
- Wave-Field-Synthesis

Evolution:

Beamforming
(Mic. array)

- problem in audio

$$\begin{aligned} f \uparrow \uparrow &\Rightarrow N \cdot \Delta x \uparrow \uparrow \\ f \uparrow \uparrow &\Rightarrow \Delta x \uparrow \uparrow \\ &\downarrow \text{separación} \\ &\text{baja} \quad \text{gran apertura} \end{aligned}$$

ICA:

corte computacional elevadísimo

CASA:

perceptual sound processing principles based on listening experiments

problem: very specific e.g. speech singing voice

Sparse methods:

Idea: baja probabilidad de que fuentes distintas sean simultáneas en tiempo o frecuencia o ambos

↳ time/frequency masking

MIXTURES

→ Audio Sources

• VOZ: secuencia de fonemas

- ↳ cuerdas vocales (pitch fundamental + armónicos)
- ↳ ruidos de susurro (también tiene formantes)
- ↳ transitorios (labios)

Male around 140 Hz
Female 200 Hz

parciales

• MUSIC: musical instruments + singers

- ↳ tonos/notas = harmonic partials → acordes = mezcla notas
- ↳ transient signals (ej: tambores)

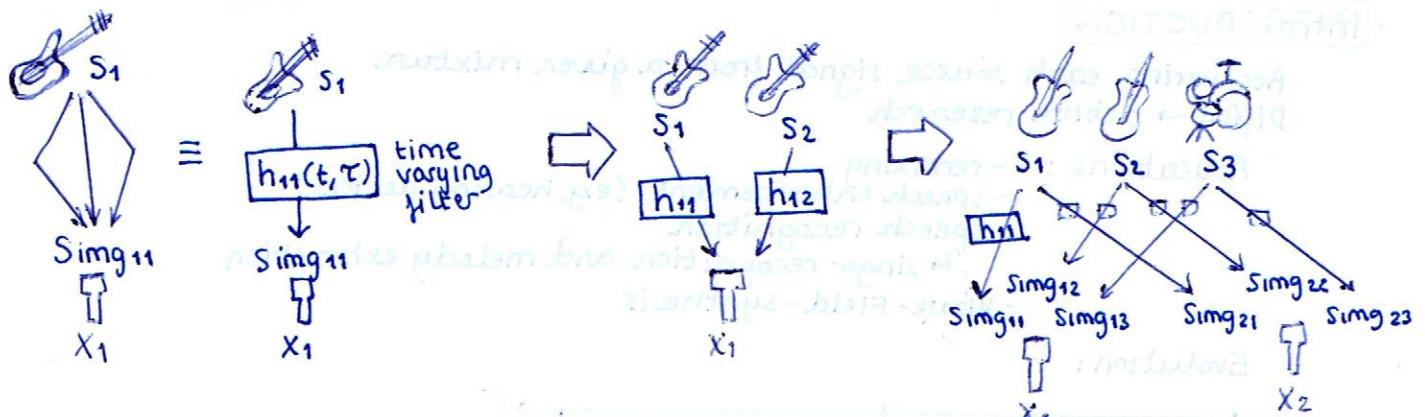
• SONIDOS NATURALES: muy distintos ej: pitito de coche
ej: lluvia

• STUDIO RECORDINGS: multitrack + downmix recording

• LIVE RECORDINGS: Puede intentarse la separación con micrófonos
ej: 2 micrófonos X-Y → apuntando ↗ ↘ cardióide en 8
ej: 2 micrófonos spaced → sencillo y rápido ↗ ↘

• SYNTHETIC MIXTURES: tabla de mezclas

The mixing process



En SSS típicamente se intenta estimar $S_{img\ ij}(t)$ efecto sala

• si trato de estimar directamente $S_j(t)$ se llama derreverberación

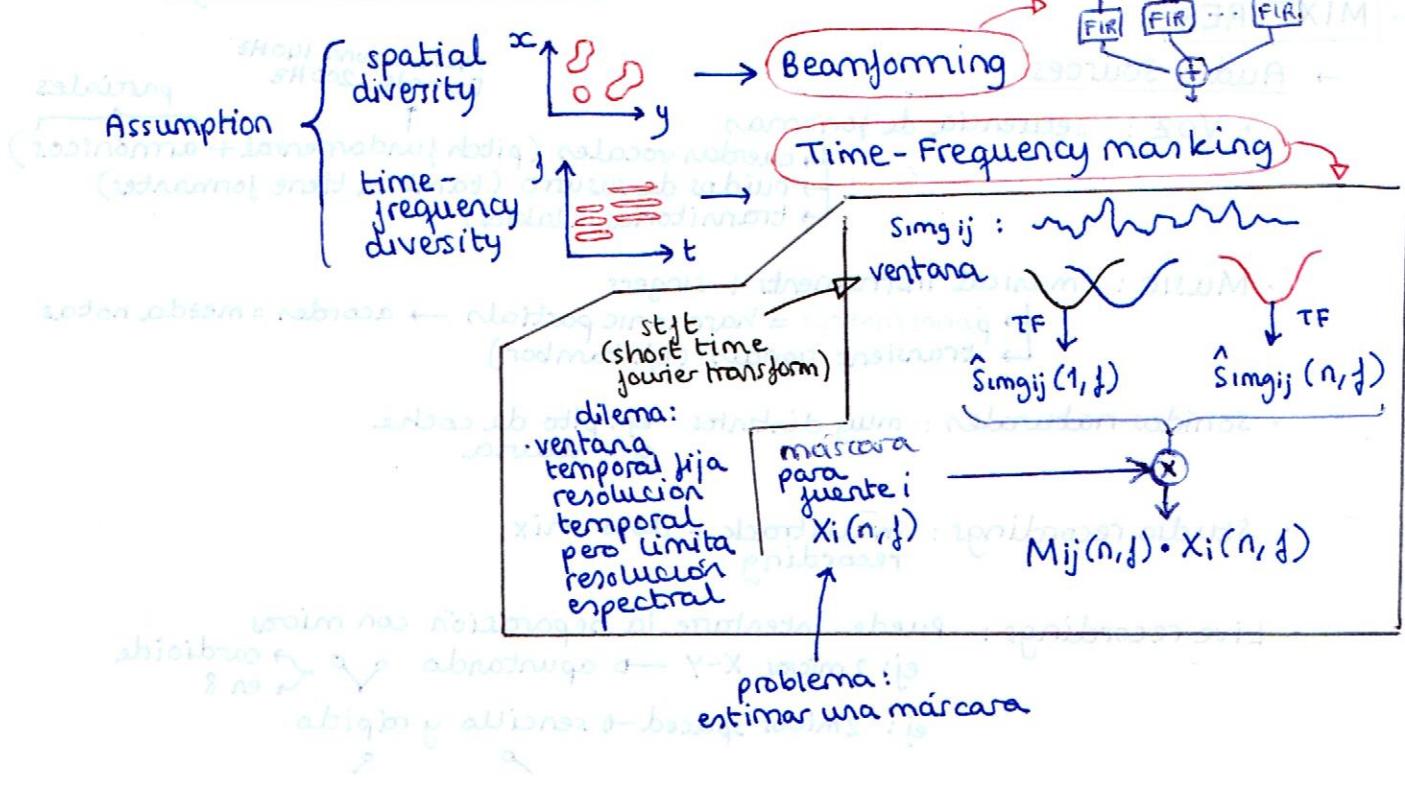
$$x_i(t) = \sum_{j=1}^J S_{img\ ij}(t)$$

$\underbrace{S_j(t) * h_{ij}(t, \tau)}$

BLIND SOURCE SEPARATION

→ separar haciendo el menor número de suposiciones posibles

FILTERING TECHNIQUES



→ MULTI-CHANNEL SEPARATION

• ICA (Independent Component Analysis)

Assumption: the sources are statistically independent

↳ minimize mutual information between estimated source signals

↳ find optimal demixing matrix (MAP)

$$\begin{bmatrix} 1 & 0 & \dots \\ 0 & 1 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \Delta = (1, 0) \text{ det}$$

• ADReSS (Azimuth Discrimination and Resynthesis)

Idea stereo:

$$\begin{aligned} \text{Left } l(t) &= \underbrace{\text{voz}}_{\text{guitarra}} + \underbrace{\text{guitarra}}_{\text{voz}} \\ \text{Right } r(t) &= \underbrace{\text{guitarra}}_{\text{voz}} + \underbrace{\text{voz}}_{\text{guitarra}} \end{aligned}$$

$$l(t) - g_j r(t) \rightarrow \cancel{\text{cancel}}$$

causes the cancellation of the j-th source in the left channel

$$r(t) - g_j l(t) \text{ causes the cancellation of the j-th source in the right channel}$$

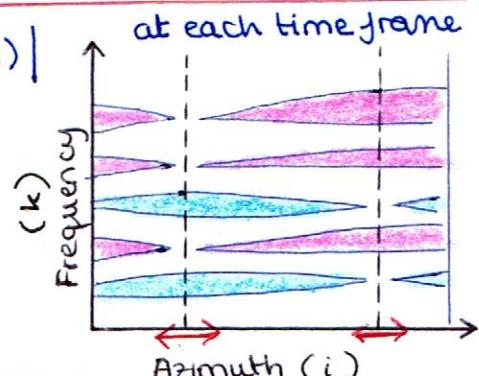
We know how to cancel a source, but how can we recover it?

①

$$A_{Zr}(k, i) = \left| L(k) - \frac{i}{\beta} R(k) \right|$$

$\begin{bmatrix} g(i) \\ i \in [0, \beta] \end{bmatrix}$

y lo mismo con canal left $A_{Zl}(k, i)$



② Convertir nulos en picos

Buscar picos en un azimuth subspace width

↳ ya sabemos a qué frecuencias está cada fuente en un determinado time frame

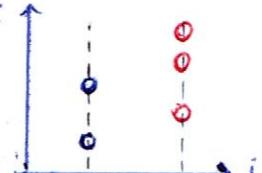
$$A_{Z'}(k, i) = \begin{cases} A_Z(k)_{\max} - A_Z(k)_{\min} & \text{if } A_Z(k, i) = A_Z(k)_{\min} \\ 0 & \text{resto} \end{cases}$$

③ señal separada:

$$Y_R(k) = \sum_i A_{Zr}'(k, i)$$

$$Y_L(k) = \sum_i A_{Zl}'(k, i)$$

modulo (la que se cogiera del original $R(k), L(k)$)



(+) Good quality of separation

(+) works with stereo → casi todo está en estéreo

(+) Very simple

Problems: - sources with same pan position are extracted together
- musical noise in extracted sources

You can exploit spatial clues

Mirar estadísticos de alto orden (> 2) (mas allá de la varianza)

Ejemplo:

- una gaussiana da cero

- suma de 2 gaussianas no da cero, pero si las separo cada una da cero

DUET (Degenerate Unmixing Estimation Technique)

Assumption: only one source is present at any time-frequency point (W-disjoint orthogonality)

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

It is possible to determine the source in each point from the spatial location information

$$\text{Interchannel Intensity Difference} \quad \text{IID}(n, f) = 20 \log \frac{X_2(n, f)}{X_1(n, f)}$$

$$\text{Interchannel Phase Difference} \quad \text{IPD}(n, f) = \Delta \left[\frac{X_2(n, f)}{X_1(n, f)} \right]$$

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

Problems:

- Rarely holds
- musical noise artifacts
- it cannot deal with convolutive mixtures where mixing delay is larger than one sample

SINGLE CHANNEL SEPARATION

- Spatial clues cannot be included
- more challenging problem

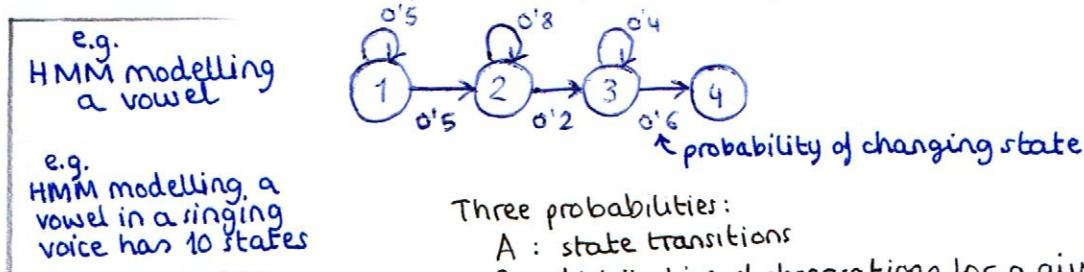
• Hidden Markov Models (HMM)

~~Markov model: states correspond to random processes whose outcomes are observable~~

HMM: the states are not observable directly but rather have a certain probability distribution of the observation.

In Speech:

- Observations are significant speech parameters
- Each state models a particular phoneme or particular chord



steps:

- ① Train individual HMM beforehand to get A , B , π
- ② calculate state path (most likely succession of states) (e.g. Viterbi)
- ③ Each state corresponds to a short-term spectrum of the source

Three probabilities:

A : state transitions

B : distribution of observations for a given state

π : probability of initial state

(+) Good separation of male and female voice

(+) Acceptable separation of singing voice

(-) Difficult to determine optimum number of states for each case

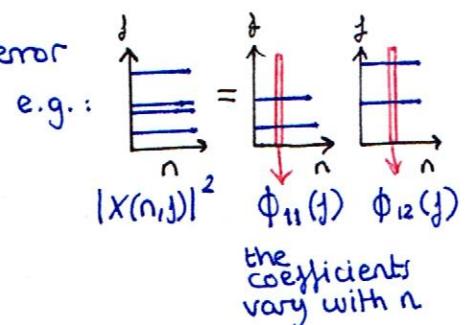
(-) Train beforehand with the solo voice

• Spectral decomposition

spectrum as weighted sum of normalised bases plus error

$$|X(n, j)|^2 = \left(\sum_j \sum_k e_{jk}(n) \cdot \phi_{jk}(j) \right) + \epsilon(n, j)$$

↓
time varying
coefficients



↳ Non-negative matrix factorization

$$X = B \cdot G$$

↑ basis
functions ↑ time varying
coefficients

• Estimate B and G
(non-negative)

• Once estimated, each source will correspond to set of basis spectra

Computational Auditory Stream Analysis (CASA)

→ Inspirado en la capacidad de los humanos para distinguir sonidos
 ↳ experimentos de audición → progresos durante 10 años

① Descomponer la muestra en componentes tiempo-frecuencia

② Agrupar los componentes en sus correspondientes fuentes

→ se realizan tres fuentes adicionales para cada grupo: MMH

se determina el orden de las fuentes

→ se determina el orden de las fuentes



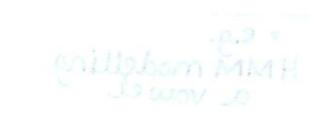
orden paralelo de particionado

orden secuencial de particionado

orden secuencial de particionado : A

orden secuencial de particionado : B

orden secuencial de particionado : C



orden secuencial de particionado : A

orden secuencial de particionado : B

orden secuencial de particionado : C

(orden A) (orden B) (orden C) (orden D)

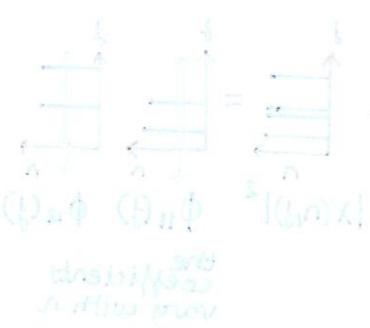
orden secuencial de particionado : D

orden secuencial de particionado : E

orden secuencial de particionado : F

orden secuencial de particionado : G

orden secuencial de particionado : H



orden secuencial de particionado : A

$$(f_{1,0})_3 + ((f_{1,0})\Phi \cdot (\text{exp} \frac{2\pi i}{T})) = \text{I}(f_{1,0})_1$$

orden secuencial de particionado : B

orden secuencial de particionado : C

orden secuencial de particionado : D

orden secuencial de particionado : E

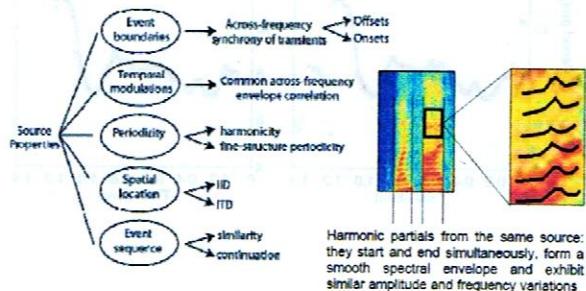
$$\frac{2}{3} \cdot \frac{8}{3} = X$$

CASA : Resumen del Proceso

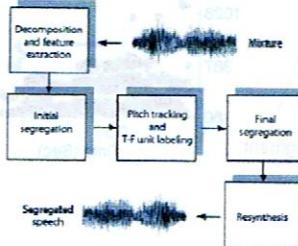
siguiente es el resumen del proceso:

10 years of progress

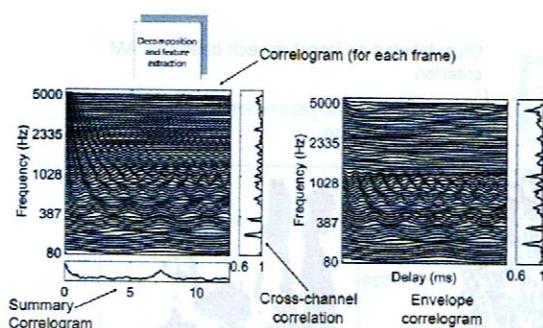
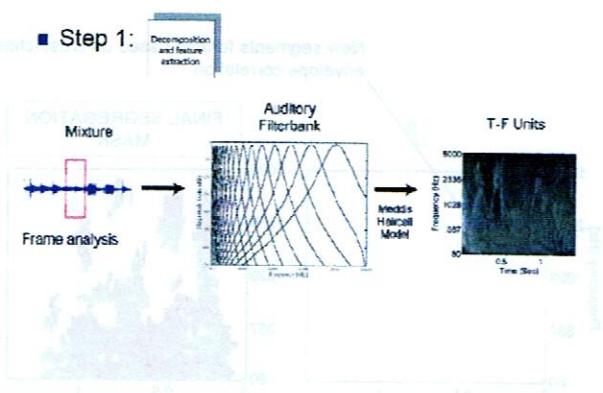
- 100 mixture set used in Cooke (1991)
- Speech + { white noise, tone, telephone, impulses, siren, music, speech, ... }
- Largely data-driven systems below
- Unrealistic corpus: mainly voiced speech + instruction



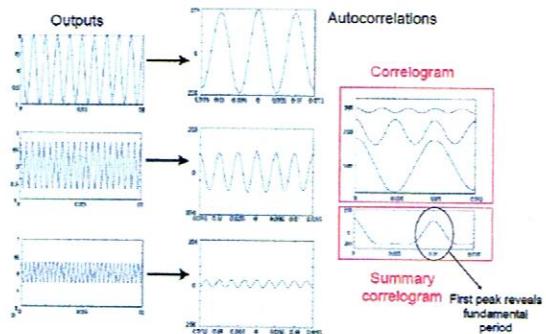
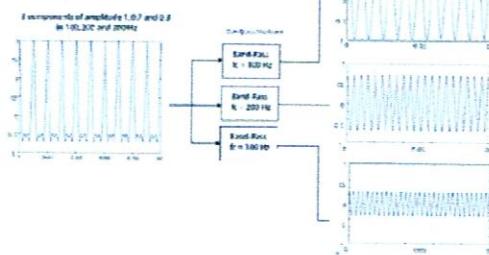
■ Monaural Speech Segregation:



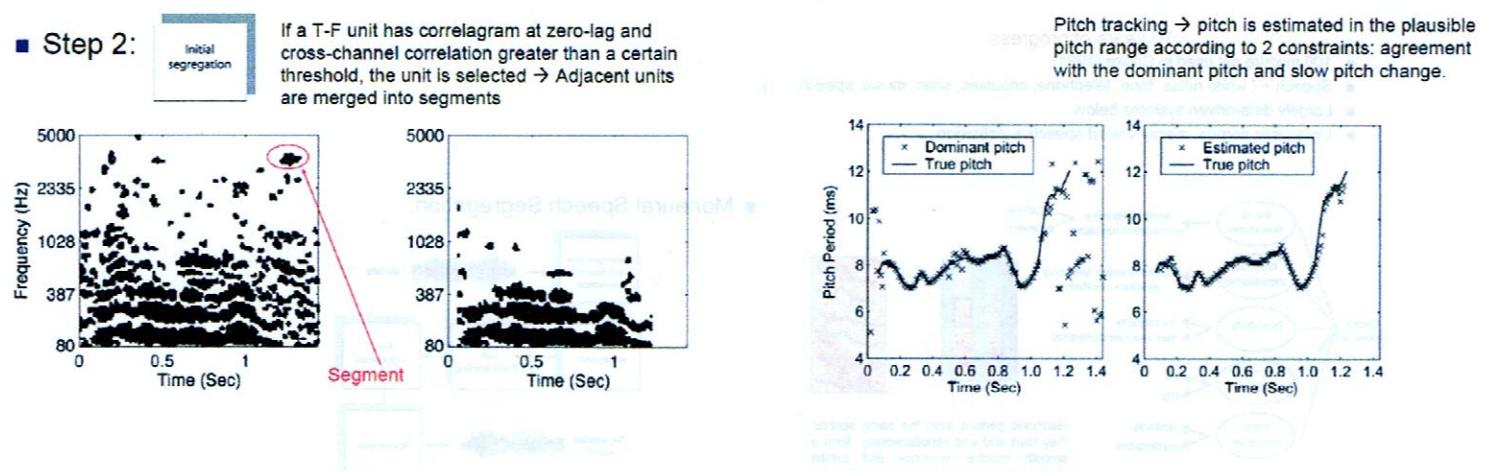
■ Step 1: Decomposition and feature extraction



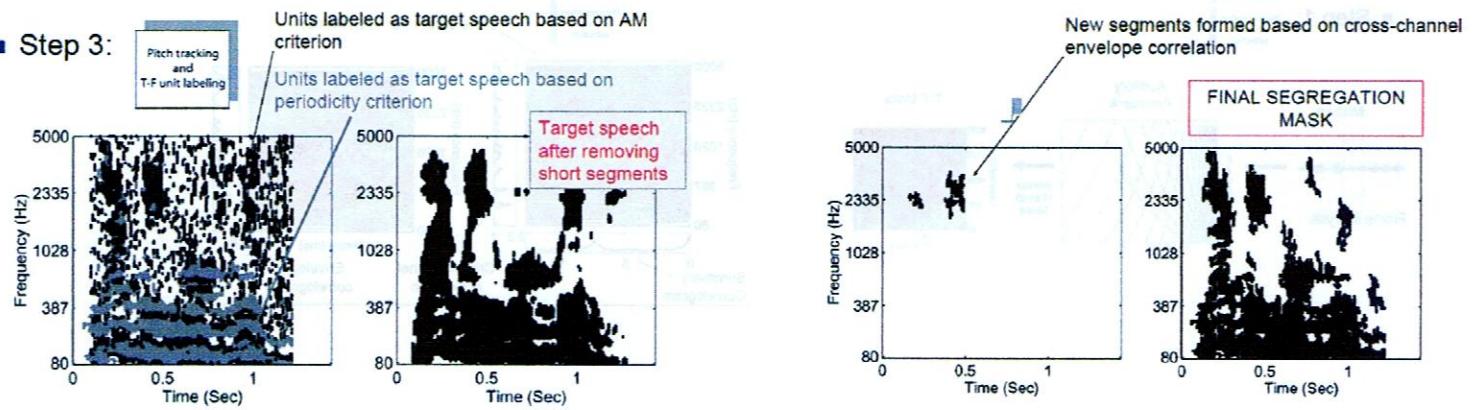
■ Correlogram:



■ Step 2:



■ Step 3:



SESIÓN 4 : SPATIAL AUDIO CODING (SAO) / SPATIAL AUDIO OBJECT CODING (SAOC)

Spatial Audio Coding (mp3 Surround)

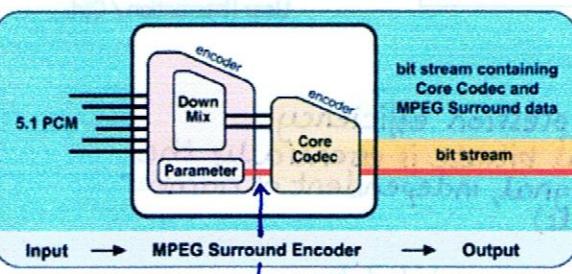
(MPEG Surround)
(MPS)

CODEC:

- Codifica multicanal en un stream stereo + info adicional

El down mix puede ser automático o externo (agusto del artista)

↑
spatial
parameters



info adicional
(reduced set of spatial parameters)

Ventajas y aplicaciones

- Tasa de bits muy cercana a estéreo → bueno para broadcast IPTV

→ Fácil de integrar por los broadcasters

→ Está en estándares de broadcast ej: DAB

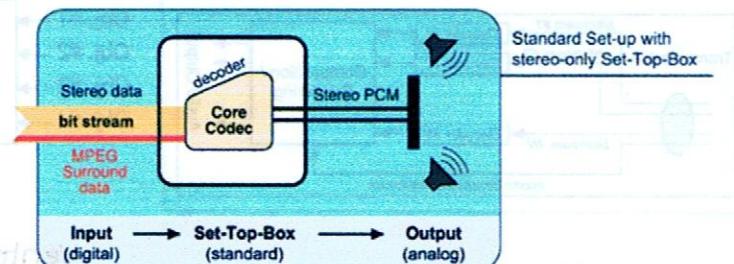
- Dispositivos móviles:
→ surround binaural

- music and video downloads
→ same music file plays on conventional stereo player

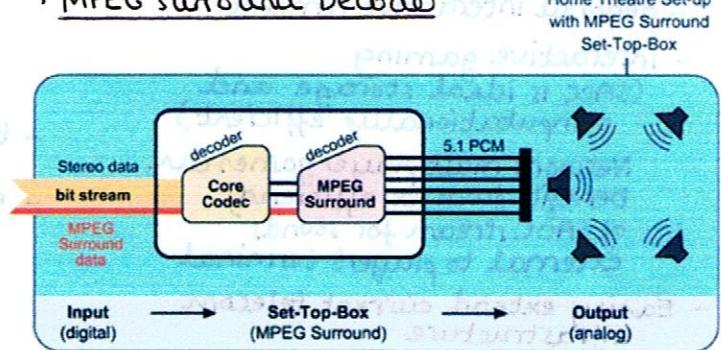
DECODEC

• Legacy decoder:

- ignora info adicional
- backward compatibility



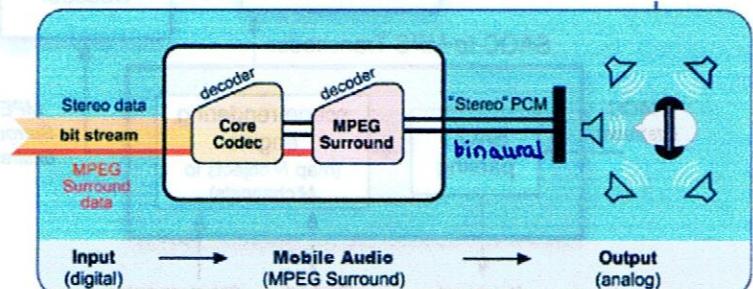
• MPEG surround Decoder



• Capacidad adicional Ensonido

- A partir del stream estéreo+parametros genera 2 canales con HRTF para usar cascos

(perfecto para dispositivo móvil!)



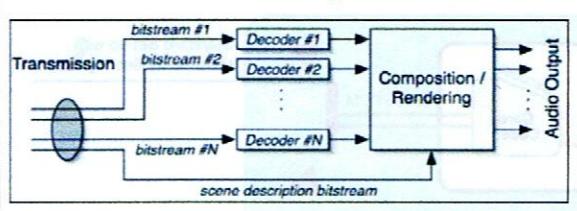
lo hace "directamente"
sin pasar por 5.1

Spatial Audio Object Coding (SAOC)

Object based representation:

- individual dry sources
- + scene description or user input

↓
Context based interactivity

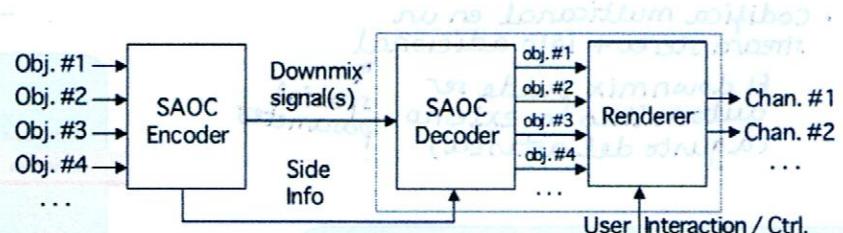


Applications:

- personal interactive remixes
- interactive gaming
(SAOC is ideal storage and computationally efficient)
- Network multiplayer games can benefit from tx efficiency of SAOC stream for sounds external to player's terminal
- Easily extend current telecom infrastructure

SAOC busca lograr la misma funcionalidad que el Object Based Representation, de forma eficiente.

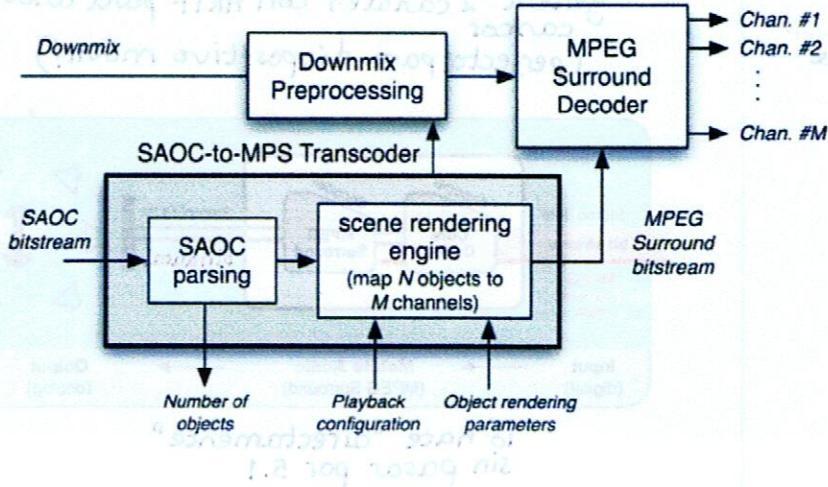
- tx :
- stereo stream
 - side info → ayuda a separar los distintos dry sources



Ventajas:

- high compression efficiency
(the total bitrate is essentially the stereo signal, independent of number of objects)
- backward compatibility
- el SAOC decoder + renderer pueden integrarse para emitir el paso intermedio de separación en objetos
→ reduce complejidad computacional

MPEG SACC Architecture



Objective:

- Integrated Decoder to avoid complexity
- N objects to M channels
 $N > M$
by simply using an M-channel MPS decoder driven by parameters