


Digital Music Representation
2) sound structure


Geoffroy Peeters
 peeters@ircam.fr
 CUIDADO I.S.T. - Ircam Analysis/Synthesis Team




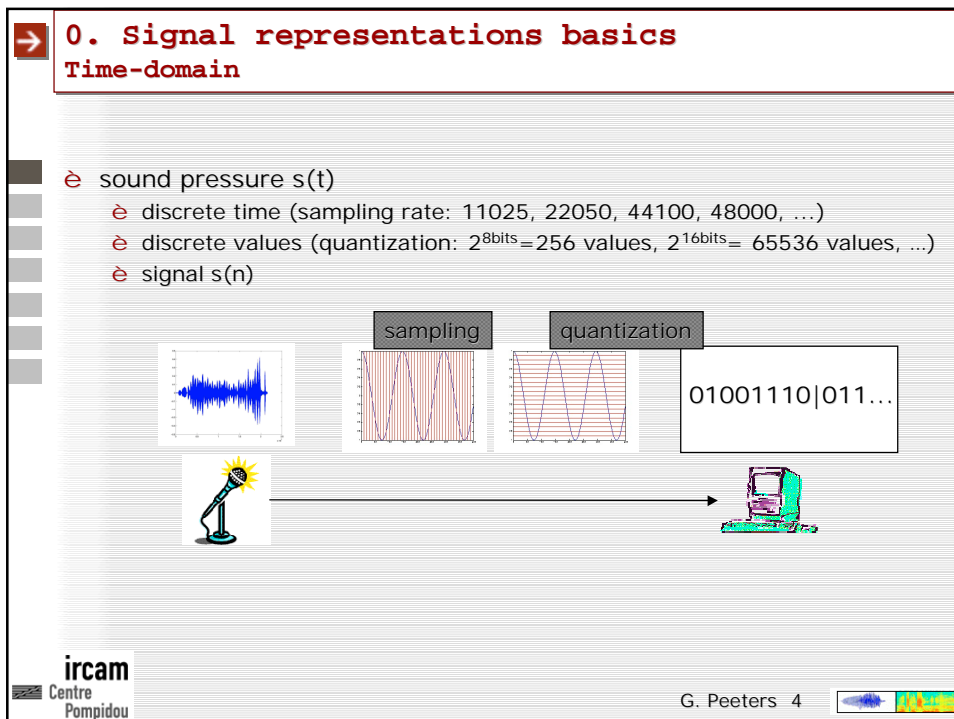
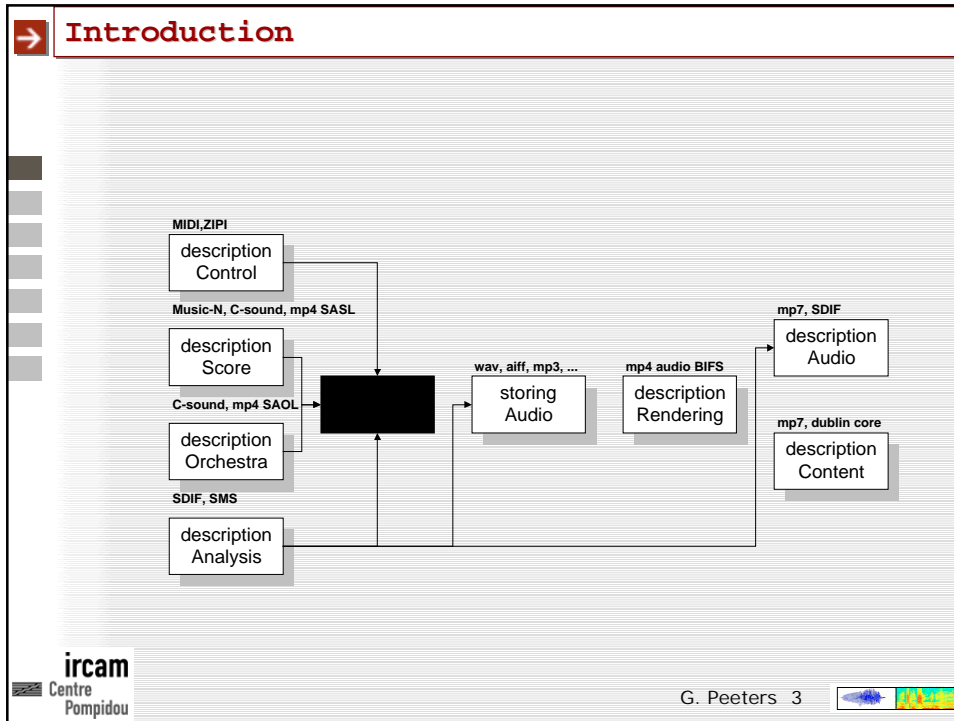
→ **Introduction**

- è Four time-scales:
 - è Eric Scheirer (1998) "The temporal structure of the pressure waveform has a significant effect on whether or not we find a sound pleasing or interesting. This structure happens on several different time scales.
 - è Notes in a musical composition change on a time scale of hundreds of milliseconds.
 - è The timbre of sounds changes on a time scale of tens of milliseconds.
 - è The actual sound waveform changes on a time scale of tens of microseconds"
 - è sound/music descriptions -> time-less
- è Representations

è to store audio:	wav, aiff, mp3, ...
è to create sound:	Music-N, C-sound, Mpeg-4 SA
è to control sound:	MIDI, ZIPI
è to render sound:	mp4 audio BIFS
è of audio analysis :	SDIF, SMS, ...
è of sound's description:	mp7, Dublin Core



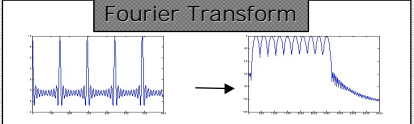
G. Peeters 2 



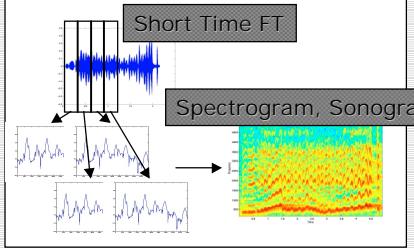
→ 0. Signal representations basics
Time-frequency/Time-scale domain

- è from the time-domain to the frequency domain
 - è most-used transforms:
 - è Fourier Transform, Fast Fourier Transform (FFT) algorithm
 - decompose the signal on a set of sinusoidal component
 - è Discrete Cosine Transform
 - è Wavelet Transform
 - decompose the signal on a set of well-localized in time and frequency functions (time-scale transforms)
- è time/frequency domain
 - è perform a transform around a given time: Short-Time Fourier-Transform
 - è 2-D representations: Time | Frequency
color represents amplitude
 - è 3-D representations: height represents amplitude

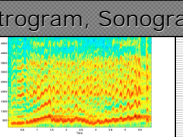
Fourier Transform



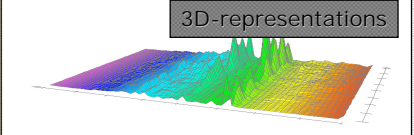
Short Time FT

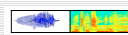


Spectrogram, Sonogram



3D-representations



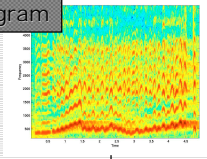
G. Peeters 5 

ircam
Centre Pompidou

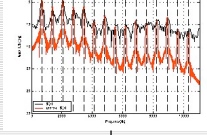
→ 0. Signal representations basics
Signal models domain

- è Most common signal models:
 - è Periodic/Harmonic model
 - è Sinusoidal additive model + noise
 - è Source / Filter decomposition
 - è LPC (linear predictive coding)
 - estimates the filter coefficients or the resonance, (in the case of speech: formant/anti-formant)
 - è CELP: coded excitation LP (source signal decomposed on a codebook)
 - è CEPSTRUM: representation of the shape of the spectrum
 - è MEL-CEPSTRUM: CEPSTRE using a Mel scale for the spectrum (perceptual scale)
 - è Granular, pitch-synchronous granular model
 - è represents the signal as a succession of grains possibly at a periodic rate

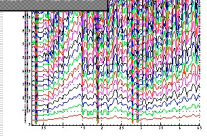
spectrogram




Sinusoidal component modeling



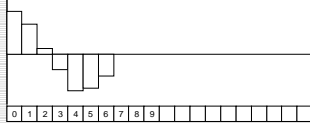
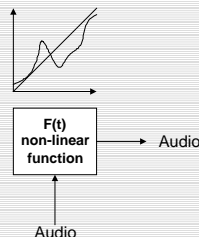
Sinusoidal tracks



G. Peeters 6 

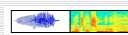
ircam
Centre Pompidou

→ 0. Signal representations basics
Synthesis algorithms

- è Wave-table synthesis
 
- è Wave-shaping synthesis
 - è apply a non-linear function to the amplitude of a sound (example: clipping that occurs when an audio amplifier is overdriven)
 - è broad range of musically useful timbres
 - è dynamic spectra

ircam
 Centre
 Pompidou

G. Peeters 7



→ 0. Signal representations basics
Synthesis algorithms

- è Additive synthesis:
 create a complex waveform by addition of basic waveforms (sinusoidal components)

spectra

waveform

OSC Freq Ampl

Analysis

Ampl envelope
 A D S R

+

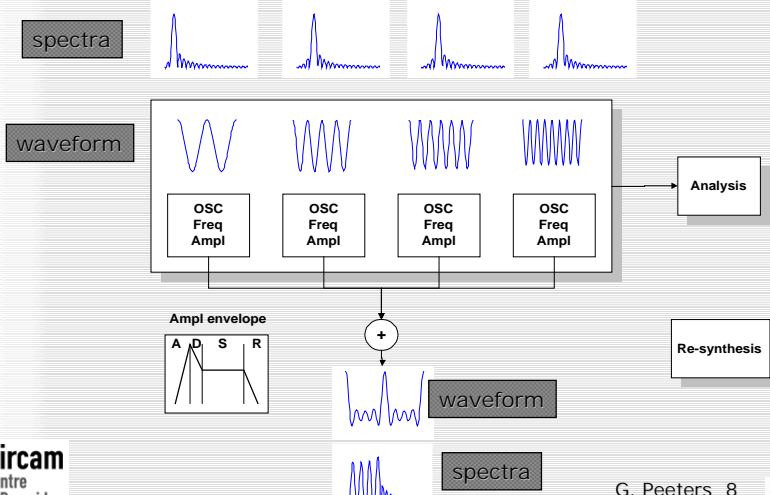
Re-synthesis

waveform

spectra

ircam
 Centre
 Pompidou

G. Peeters 8



→ 0. Signal representations basics
Synthesis algorithms

è Subtractive synthesis:
 remove components by filtering a spectrally rich waveform
 (waveform=white noise, saw-tooth, triangular, square waveform)

Ampl envelope
 A D S R

ircam
 Centre
 Pompidou

G. Peeters 9

→ 0. Signal representations basics
Synthesis algorithms

è Source/Filter synthesis:
 based on physics: a periodic "source" signal is filtered by a filter
 (ARMA filter) (see vowel synthesis = formant synthesis)

Source Signal

ARMA filter

Audio

ircam
 Centre
 Pompidou

G. Peeters 10

0. Signal representations basics
Synthesis algorithms

FM synthesis:

- Principle:
 - low frequency oscillator modulating a VCO -> vibrato effect
 - high frequency oscillator modulating a VCO -> complex spectrum (see Bessel functions)
 - Complex spectrum: depends on the magnitude of the modulation
- a set of "operators"
- Each operator consists of:
 - A digitally controlled oscillator (DCO)
 - An amplifier
 - An envelope generator
- minimum: two operators
 - a modulator
 - a carrier

The diagram illustrates the FM synthesis process. It features two main oscillator blocks: 'OSC MOD' (modulator) and 'OSC CAR' (carrier). The 'OSC MOD' block receives 'amplitude of modulation' and 'modulating frequency' as inputs. The 'OSC CAR' block receives 'carrier frequency' and 'amplitude of carrier' as inputs. The outputs of both oscillators are combined at a summing junction (+). The resulting signal is labeled 'Audio'. A 'spectra' plot is shown to the right, displaying the complex spectrum of the synthesized signal. The Ircam logo and 'Centre Pompidou' are visible in the bottom left, and 'G. Peeters 11' is in the bottom right.

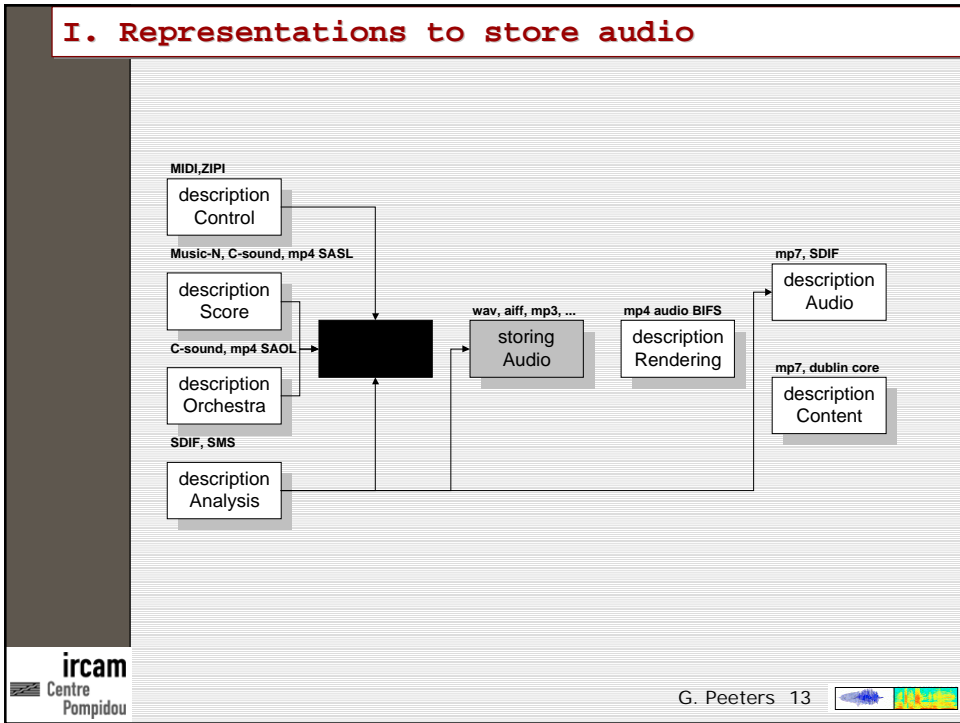
0. Signal representations basics
Synthesis algorithms

Physical modeling synthesis:

- represent the object producing the sound instead of the produced sound itself
- Most well-known physical model: Karplus-strong (plucked string)
- require a different model for each instrument
- Various modeling techniques for physical models
 - mass-spring paradigm
 - modal synthesis
 - wave-guide synthesis

playing

The diagram shows a 'Wave-guide model of clarinet-like instrument'. It starts with 'mouth pressure' entering a 'Non-linear excitation' block. The output of this block is split into two paths. One path goes to a 'Delay line z⁻¹' block, which then feeds into an 'Attenuation Filter'. The other path goes to another 'Delay line z⁻¹' block, which also feeds into an 'Attenuation Filter'. The outputs of both 'Attenuation Filter' blocks are combined and feed into a 'Radiation Filter', which produces the final 'audio' output. The Ircam logo and 'Centre Pompidou' are visible in the bottom left, and 'G. Peeters 12' is in the bottom right.



→ I. Representations to store audio

- è Non-compressed audio representations
 - è Representation of the sound (sampled and quantified) pressure: $s(n)$
 - è Using PCM (Pulse Code Modulation)
 - è Store the audio signal and related information for reading it
 - è Most format are based on "chunks"
 - è data chunks
 - è information chunks
= information needed in order to read the file: big/little endian, sampling rate, quantization, nbchannels
 - è extra chunks:
 - marker chunks
 - instrument chunks (play mode, loop points)
 - midi data chunk,
 - comments chunk (marker+text)
 - application specific chunks
 - meta-information: name, author, copyright, annotation chunks
 - wav: bitmap
 - è Examples:
 - aiff (Audio Interchange File Format) (Mac)
 - wav PCM (RIFF wav) (Windows)
 - au, snd (sun/next station)

G. Peeters 14

I. Representations to store audio

AIFF structure

```

FORM AIFF Chunk
ckID = 'FORM'
formType = 'AIFF'
    Common Chunk
    ckID = 'COMM'
    Sound Data Chunk
    ckID = 'SSND'
        
```

Chunk format

```

ckID
ckSize } header info
data
        
```

Common chunk


```

#define CommonID 'COMM' /* ckID for Common Chunk */
typedef struct {
    ID ckID;
    long ckSize;
    short numChannels;
    unsigned long numSampleFrames;
    short sampleSize;
    extended sampleRate;
} CommonChunk;
        
```

Example

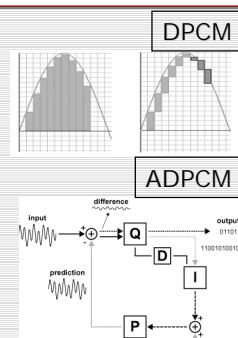
```

FORM AIFF
ckID 'FORM'
ckSize 176526
formType 'AIFF'
Common Chunk
ckID 'COMM'
ckSize 18
numChannels 2
numSampleFrames 88200
sampleSize 16
sampleRate 44100.00
Marker Chunk
ckID 'MARK'
ckSize 34
numMarkers 2
id 1
position 44100
markerName 'B | b | e | g | i | n | o | f | f | o'
id 2
position 88200
markerName 'B | e | n | d | i | n | g | o | f | f | o'
Instrument Chunk
ckID 'INST'
ckSize 20
baseNote 60
detune -3
lowNote 57
highNote 63
lowVelocity 1
highVelocity 127
gain 6
sustainLoop.playMode 1
sustainLoop.beginLoop 1
sustainLoop.endLoop 2
releaseLoop.playMode 0
releaseLoop.beginLoop -
releaseLoop.endLoop -
Sound Data Chunk
ckID 'SSND'
ckSize 176408
offset 0
blockSize 0
soundData ch 1 ch 2 ... ch 1 ch 2
first sample frame BB200th sample frame
        
```

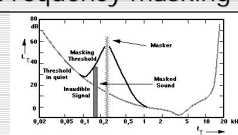


I. Representations to store audio

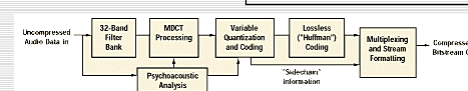
- è Compressed audio representations
 - è Not a direct representation of the sound pressure
 - è Examples:
 - è wav mu-law, a-law, DPCM, IMA ADPCM
 - è aiff-c, aifc
 - è MPEG-1 Layer III (subband-coding, masking)
 - è MPEG-2 AAC
 - è QT (QuickTime)
 - è WMA (Windows Media Audio)
 - è RA (Real Audio)
 - è SWA (Shockwave Audio)
 - è ATRAC3 (Adaptive Transform Acoustic Coding 3)




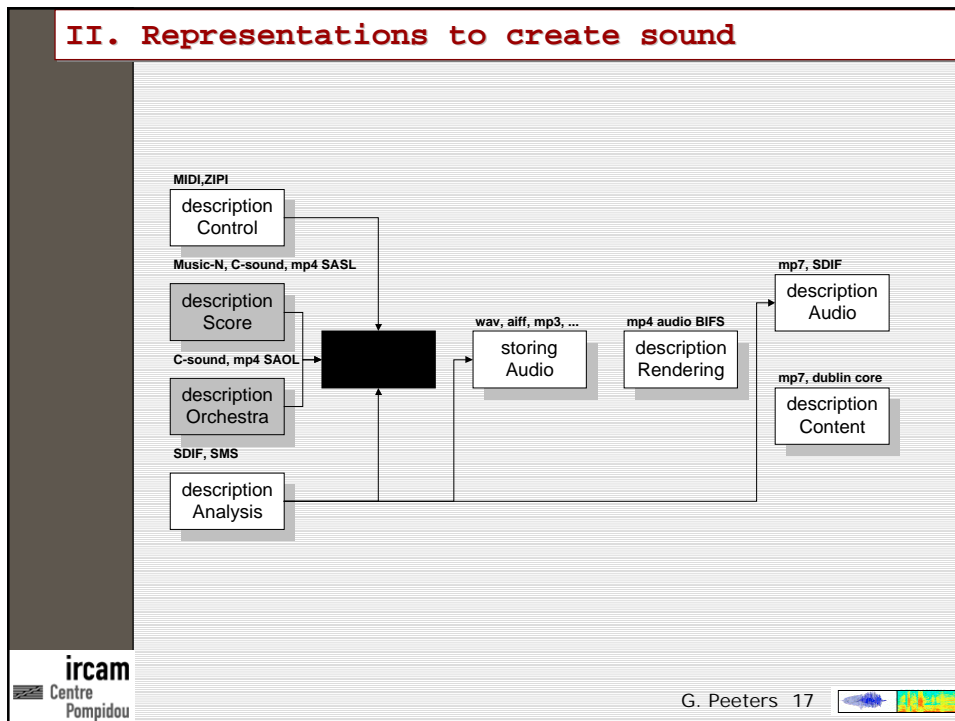
Frequency masking



Perceptual coders







→ II. Representations to create sound

C-Sound

- è History of computer sound synthesis
 - è 1957: Max Mathews: Music I (IBM 704 computer)
 - è ...
 - è 1960: Music III: concept of modular instruments (unit generators)
 - è 1962: Music IV: orchestra/score model
 - è 1969: Music V
 - è 1980: Richard Moore: C-music
 - è 1985: Barry Vercoe: C-sound
- è Concept:
 - è Richard Boulanger
 - è "C-sound is a sound renderer. It works by first translating a set of text-based instruments, found in the orchestra file, into a computer data-structure that is machine-resident. Then, it performs these user-defined instruments by interpreting a list of note events and parameter data that the program "reads" from: a text-based score file, a sequencer-generated MIDI file, a real-time MIDI controller, real-time audio, or a non-MIDI devices such as the ASCII keyboard and mouse."

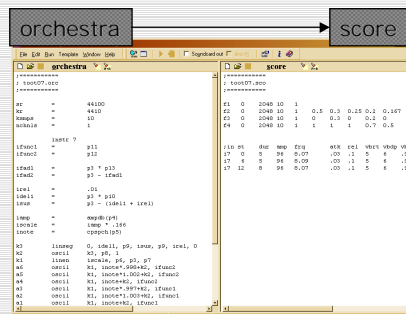
G. Peeters 18

→ II. Representations to create sound
C-Sound

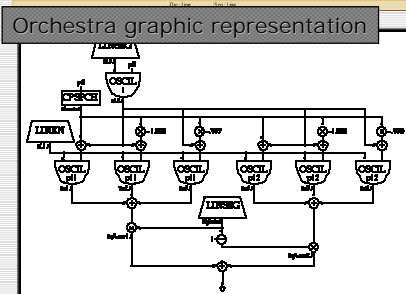
- è C-sound concepts:
 - è distinction between orchestra and score languages
 - è orchestra: modular instruments, "unit generator", macro language
 - è distinction between sample rate and control rate
 - è i-rate variables: changed at the note rate
 - è k-rate variables: changed at the control signal rate
 - è a-rate variables: changed at the audio signal rate
 - è compilation creates algorithm which then creates sound
 - è latest versions: MIDI compatible, real-time use
 - è Syntax example:
 - è Unit generator list:
 - è sinusoidal generator
 - è reverb generator, ...
 - è ...

orchestra

score



Orchestra graphic representation

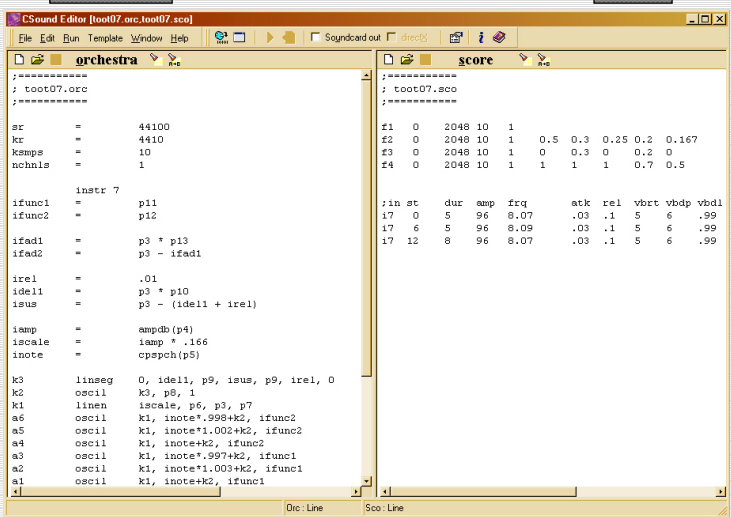


ircam
Centre
Pompidou


→ II. Representations to create sound
C-Sound

orchestra

score


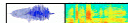


ircam
Centre
Pompidou

G. Peeters 20 

→ **II. Representations to create sound**
MPEG-4 Structure Audio (Scheirer, 1998)

- è History:
 - è "Netsound" (Casey, Smaragdis 1996): using C-sound to transmit sounds over networks
 - è C-sound intellectual properties issues -> Scheirer develops SA
- è Structure Audio: "a Music Synthesis Language"
 - è SAOL: Structure Audio Orchestra Language
 - è SASL: Structure Audio Score Language
- è "highly readable", "highly modular", "highly expressive", "highly functional"
- è more understandable and concise than C-sound ("C-like" language)
- è new features
 - è 100 built-in "unit generator", may be extended with new "unit-generator"
 - è "cpuload" -> allows Dynamic voice-stealing algorithms in the orchestra
 - è ...
- è SAOL may be controlled by
 - è SASL
 - è MIDI-file
 - è real-time MIDI events

 G. Peeters 21 

→ **II. Representations to create sound**
MPEG-4 Structure Audio (Scheirer, 1998)

SAOL

```

instr beep3(pitch) {
  importe ksig amp, off;           // controllers
  ksig vol, init;
  table wave(harm,2048,1);
  asig sound;

  if (!init) {                    // first time we're called
    amp = 0.5; init = 1;
  }
  if (off) { turnoff; }           // we got the 'off' control
  vol = port(amp,0.2);           // smooth the volume signal
  sound = oscil(wave,pitch);
  output(sound * vol);
}

```

SASL



```

n1: 0.0 beep3 1 440               // first note
    0.5 beep3 1 480               // second note

n2: 1.0 beep3 -1 220              // third note - unbounded initial duration
n2: 1.0 beep3 -1 440              // fourth note
n3: 1.0 beep3 -1 660              // fifth note

2.0 control n2 amp                // applies to third and fourth notes
2.5 control n2 amp 0.5
3.0 control n3 amp 0.2            // applies to fifth note
3.0 control n2 amp 0.2
4.0 control n2 off 1
4.0 control n3 off 1

```

 G. Peeters 22 

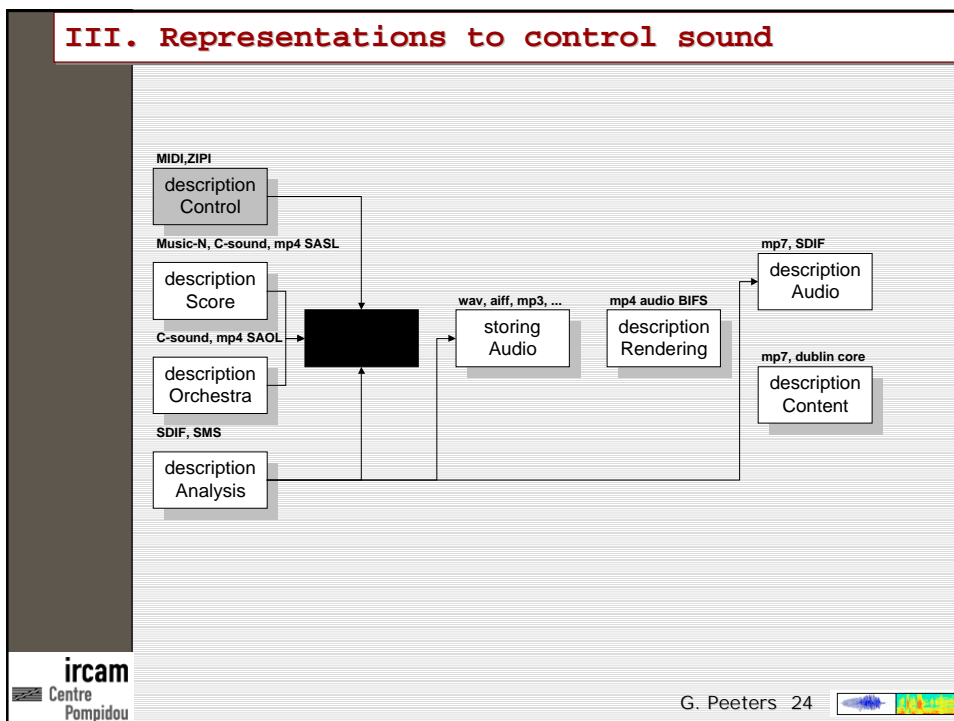
→ II. Representations to create sound
"Unit generator", Modular graphical conception Synth

è Nyquist, CLM, SuperCollider, jMax, MAX-MSP, Pure Data, Reaktor, ...

The image displays three software interfaces used for sound creation. On the left is **jMax**, showing a 'PAF vowel synthesis' window with various sliders and buttons for generating sounds. In the center is **Max/MSP**, showing a complex modular synthesis patch with numerous interconnected objects. On the right is **CsoundMax**, a real-time GUI for Csound in Max/MSP, featuring sections for 'Sound Performance', 'P-Field Control', 'Keyboard Control', and 'Audio Performance'.

ircam
Centre
Pompidou

G. Peeters 23



→

III. Representations to control sound

MIDI

- è MIDI (Musical Instrument Digital Interface) 1982/1983
 - è Goal: allow musicians to connect synthesizers together
 - è How: instruct a synthesizer which sounds to use, which notes to play, how loud to play it, ...
 - è Characteristic: unidirectional asynchronous bit stream at 31.25 Kbits/sec
 - è Today's usage: used in digitized audio in games and multimedia applications
 - è Advantages
 - è storage place: transmit only the playing information not the sound (1 minutes of audio: 10Mbytes in midi only 10Kbytes)
 - è easiness of manipulation

G. Peeters 25

→

III. Representations to control sound

MIDI

- è MIDI messages
 - è Channel Voice Messages (16 channels) carry musical performance data
 - è examples: Note On, Note Off, Polyphonic Key Pressure, Channel Pressure, Pitch Bend Change, Program Change, and the Control Change messages
 - è Mode Messages affect the way a receiving instrument will respond to the Channel Voice messages
 - è examples: reset all controller, Omni On/Off
 - è System messages
 - è System Common Messages: include MTC Quarter Frame, Song Select, Song Position Pointer, Tune Request, and End Of Exclusive (EOX) -> Midi Time Code
 - è System Real Time Messages: used to synchronize all of the MIDI clock-based
 - Timing Clock, Start, Continue, Stop, Active Sensing, and the System Reset message
 - è Midi limitations:
 - è 1) limited bandwidth -> synchronization
MIDI = 31.25 Kbit/s, serial transmission (possible that musical events which originally occurred at the same time may not actually be played at exactly the same time)
 - è 2) dedicated to keyboard-like instruments

G. Peeters 26

➔

III. Representations to control sound

Midi Files - GM - GS - XG

- è Standard Midi Files:
 - è why ? midi has been developed for real-time, need to add a "time-stamping" for the MIDI messages when stored
 - è Format 0: single track, Format 1: several tracks, Format 2: several independent patterns
- è General Midi (GM):
 - è why ? problem with "pure" MIDI: no standard for the relationship between patch numbers and specific sounds for synthesizers
 - è definition of a General MIDI Sound Set (a patch map): 128 basic patches, MIDI Channels 1-9 and 11-16
 - è definition of a General MIDI Percussion map (mapping of percussion sounds to note numbers): MIDI Channel 10
- è Roland GS
 - è why ? extension of GM because of GM limitations
 - è additional controllers (synthesizer parameters, effect parameters)
 - è scheme of variation tones (banks of GM sounds)
 - è scheme of tone 'fallback' (for non-existent banks)
- è Yamaha XG
 - è a more powerful extension of GM

G. Peeters 27

➔

III. Representations to control sound

ZIPI

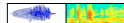
- è ZIPI (Gibson/ Zeta Music, CNMAT)
- è Motivation:
 - è MIDI is based on keyboard controlling: a note cannot start without a pitch
 - è not appropriate for controlling from guitar, violin, sustained instruments
- è ZIPI characteristics
 - è network oriented, dynamic assignment
 - è larger bandwidth: lower bandwidth=10 times MIDI
- è ZIPI
 - è Music Parameter Description Language (MPDL)
 - è MIDI
 - è Data dump
 - è Digital Audio

G. Peeters 28

→ III. Representations to control sound
ZIPI

- è ZIPI Music Parameter Description Language (MPDL)
 - è 1) note address
 - è MIDI: address=channel, or address=pitch's note
 - è ZIPI MPDL: address=note whatever the pitch
 - è ZIPI MPDL: group of group of notes:
 - “notes” (127) in “instruments” (127) in “families” (63) = 1.016.127 addresses
 - MIDI: no way to articulate an entire chord with a single message (note or channel based message)
 - è each ZIPI device has its own address space
 - è 2) note descriptors
 - è syntax: descriptor ID | value
 - è example of descriptor:
 - articulation, pitch, frequency in Hz, Amplitude, Loudness, Brightness, Even/odd harmonic balance, roughness, spatialization azimuth angle, timbre space X dimension, ...
 - articulation: note sounds, note re-attacks
 - program change
 - Higher order-messages:
 - modulation, housekeeping (defines notes priority), querying a synthesizer, comments, time-tags (ensure synchronization when streaming ZIPI (minimum delay to respect))
 - è 3) in ZIPI: distinction between Controller and Synthesizer Messages (in MIDI both are mixed)

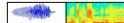
ircam
 Centre
 Pompidou

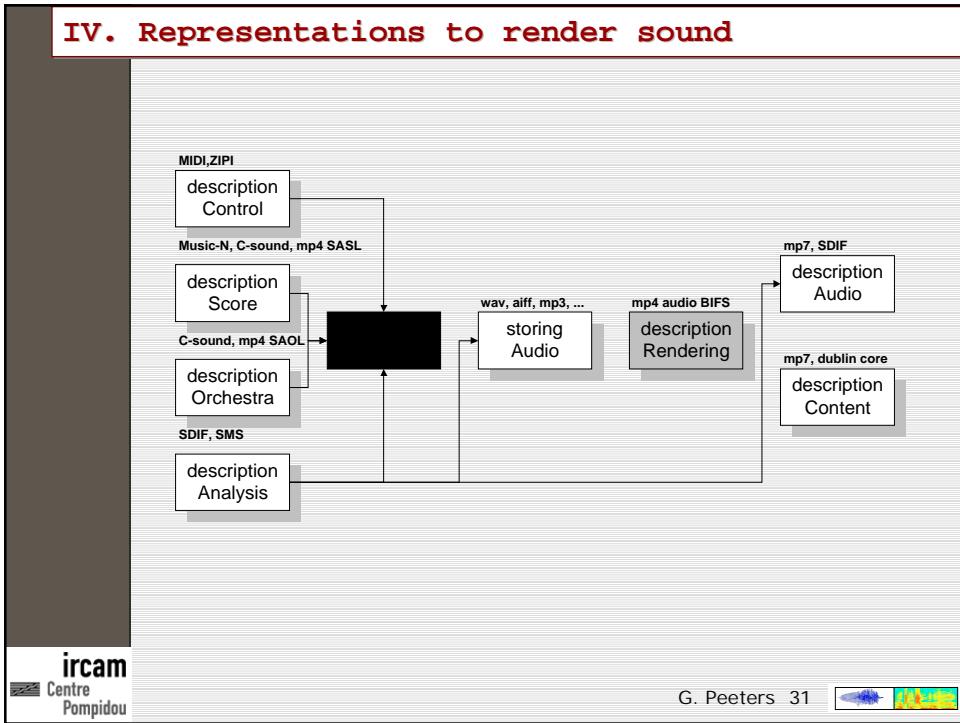
G. Peeters 29 

→ III. Hybrid sound/control representations
MOD/RMF

- è MOD
 - è tracker file format
 - è contains both
 - è musical score information
 - è actual instrument sound samples that are used to play a tracker song.
- è RMF (Rich Music Format)
 - è three types of information:
 - è highly compressed sound samples,
 - è control data (playback instructions, like MIDI),
 - è specifications for interactivity

ircam
 Centre
 Pompidou

G. Peeters 30 



→ IV. Representations to render sound

- è Multi-channel file formats
 - è Dolby Digital (Dolby AC-3) 5.1
 - è 5 discrete (independent) channels (center, left, right, surround left, surround right) (from 20Hz to 20,000 Hz)
 - è 6th channel dedicated for low frequency effects (LFE) (3 Hz to 120 Hz).
 - è DTS Digital Surround 5.1
 - è ...
- è Rendering description formats
 - è MPEG-4 Audio BIFS
 - è BIFS (Binary Format for Scenes) in SNHC (Synthetic/Natural Hybrid Coding) based on VRML
 - è Sound attaches to objects: sound node used to put sound in the scene
 - è MPEG-4 Version 2: Advanced Audio BIFS
 - è Sound node
 - è Sound BIFS
 - Acoustic/physical approach
 - Perceptual approach (SPAT)

Acousmonium

Espace de projection

ircam
Centre
Pompidou

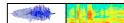
G. Peeters 32

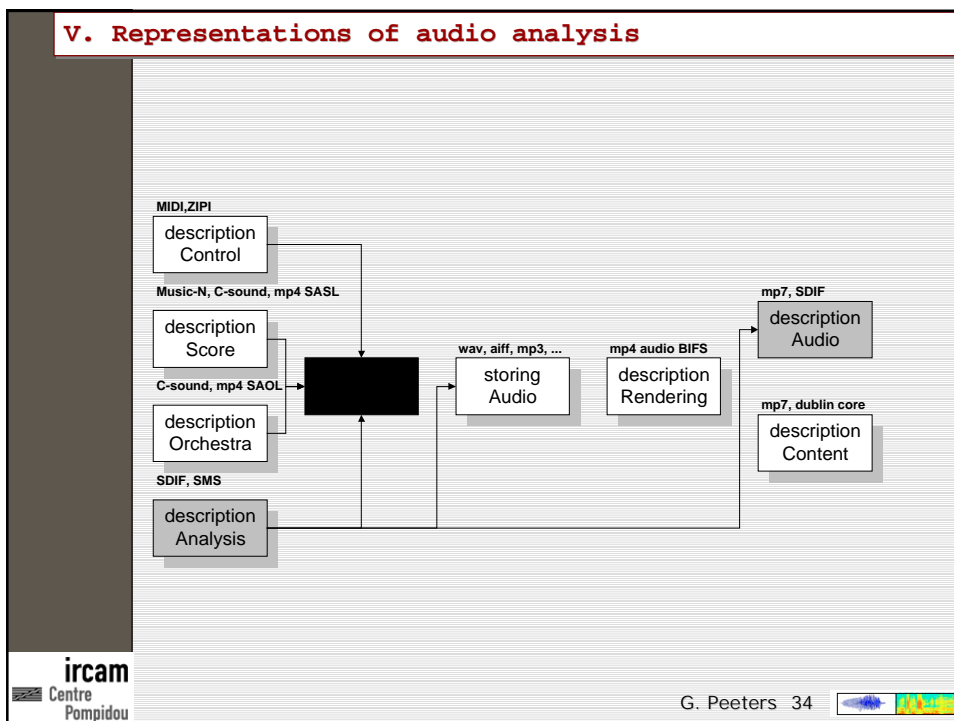
→ IV. Representations to render sound

è Two approaches:

<p>è Acoustic/physical approach</p> <ul style="list-style-type: none"> è Directive sound <ul style="list-style-type: none"> è source è location, direction, intensity, speedOfSound è angles, directivity filters è useAirabs, spatialize, roomEffects è Acoustic Scene <ul style="list-style-type: none"> è late reverberation, reverbTime, reverbLevel, reverbDelay è center, size of bounding box è AcousticMaterial <ul style="list-style-type: none"> è acoustic reflectivity and transmissions filters è ... 	<p>è Perceptual approach (SPAT)</p> <ul style="list-style-type: none"> è Directive Sound è Perceptual Parameters <ul style="list-style-type: none"> è source presence, è warmth, è brilliance, è room presence, è envelopment, è running reverberance, è late reverberance, è heaviness, liveness
--	--

ircam
Centre Pompidou

G. Peeters 33 



V. Representations of audio analysis

What are the results of an audio analysis used for ?

- Parameters for re-synthesis
 - STFT
 - Sinusoidal models
 - Source/filter
 - ...

ircam Centre Pompidou

G. Peeters 35

V. Representations of audio analysis

What are the results of an audio analysis used for ?

- Parameters for re-synthesis
 - STFT
 - Sinusoidal models
 - Source/filter
 - ...
- Parameters for description
 - Fundamental frequency
 - pitch
 - Spectral centroid
 - brightness,
 - Spectral Flatness
 - Audio ID / Fingerprint
 - classification, indexing
 - high-level synthesis
- What about manual encoding ?

ircam Centre Pompidou

G. Peeters 36

V. Representations of audio analysis

- è What are the results of an audio analysis used for ?
 - è Parameters for re-synthesis
 - è STFT
 - è Sinusoidal models
 - è Source/filter
 - è ...
 - è Parameters for description
 - è Fundamental frequency
 - pitch
 - è Spectral centroid
 - brightness,
 - è Spectral Flatness
 - Audio ID / Fingerprint
 - è classification, indexing
 - è high-level synthesis
 - è What about manual encoding ?

ircam
Centre
Pompidou

Peeters 37

V. Representations of audio analysis
Graphical representations: spectral representations

V. Representations of audio analysis
Graphical representations: descriptors representations

ircam
Centre
Pompidou

V. Representations of audio analysis
SDIF

- è SDIF (Sound Description Interchange Format) 1998 (CNMAT, Ircam, IUA/UPF)
- è Meta-format (// XML)
- è Definition of a file formats and structures: chunk-based
- è Frame-based (self defined)
 - è Stream ID
 - è Frame type/matrix type examples:
 - 1FQ0 Fundamental Frequency Estimates,
 - 1STF Discrete Short-Term Fourier Transform,
 - 1TRC Sinusoidal Tracks,
 - 1RES Resonance / Exponentially Decaying Sinusoids
 - è Extensible !
- è Tools:
 - è compliant software for coding/ decoding, visualization
 - è libraries for SDIF files manipulation
- è Conclusion:
 - è allows any type of descriptions although only signal-based analysis/features are defined
- è SDIF extensions:
 - è XML DTD for Frame/Matrix definition
 - è SDIF / XML: for streams relationship
 - è SDIF to MPEG-4 cross-coder source code

SDIF structure

```

Frame
├── Frame Header
├── Matrix
│   └── Structure simple 1
│       └── ...
└── Matrix
    └── Structure simple N
                
```

SDIF viewer (IU/UPF)

ircam
Centre
Pompidou

G. Peeters

V. Representations of audio analysis

SMS

s m

s

è SMS: File Format

- è based on SDIF
- è + a set of chunks specific to SMS
- è different levels of spectral abstraction useful in a variety of applications that go from high quality audio coding to synthesis applications that use banks of spectral data
 - è SMS Generic chunk: mono, stereo, quadraphonic
 - è SMS Generic Track chunk: Track_ID, Magnitude_Threshold, Spectral_Range
 - è SMS Note chunk: Default_Magnitude_Threshold, Note_Type (pitched, unpitched), pitch
 - è SMS Note Track chunk: Note_Track_Type (whole note, attack, steady state, release, articulation)

SMS file format structure

è SMS: Score Format

- è Time statements,
- è Event statements,
- è Synth parameters
- è Controller parameters:
 Key Velocity Wind pressure, Key Number Embouchure, Key Pressure Lip pressure, Pitch-bend wheel Lip, frequency, Mod wheel Wind keypads, Switch pedal, Continuous pedal, Drum head striking X position, Drum, head striking Y position, ...

G. Peeters 41

VI. Representations of sound's description

```

    graph LR
        subgraph Description
            D1[description Control]
            D2[description Score]
            D3[description Orchestra]
            D4[description Analysis]
        end
        subgraph Storing_Audio [storing Audio]
            SA[wav, aiff, mp3, ...]
        end
        subgraph Rendering [description Rendering]
            R[mp4 audio BIFS]
        end
        subgraph Audio [description Audio]
            A1[mp7, SDIF]
        end
        subgraph Content [description Content]
            C[mp7, dublin core]
        end
        D1 --> SA
        D2 --> SA
        D3 --> SA
        D4 --> SA
        SA --> R
        R --> A1
        R --> C
    
```

G. Peeters 42

VI. Representations of sound's description

- è Meta-data description of a sound ?
 - è Content description ?
 - è Data related to the sound content or describing the sound content
 - è Data contained in the sound content or not contained in
- è Examples
 - è Information chunk in wav/aiff
 - è ID3v1/v2
 - è Dublin Core
 - è MPEG-7 Multimedia Description Interface

ircam
Centre
Pompidou

VI. Representations of sound's description
Dublin Core

- è Dublin Core (1995):
 - è description of electronic media
 - è 15 elements
 - è web resources oriented
 - è simpler than MARC

DC Element	Definition
1	Title A name given to the resource
2	Creator An entity primarily responsible for making the content of the resource
3	Subject The topic of the content of the resource
4	Description An account of the content of the resource
5	Publisher An entity responsible for making the resource available
6	Contributor An entity responsible for making contributions to the content of the resource
7	Date A date associated with an event in the life cycle of the resource
8	Type The nature or genre of the content of the resource
9	Format File format or mime type (MPEG-1, QuickTime, RealVideo...)
10	Identifier An unambiguous reference to the resource within a given context
11	Source A Reference to a resource from which the present resource is derived
12	Language A language of the intellectual content of the resource
13	Relation A reference to a related resource
14	Coverage The extent or scope of the content of the resource
15	Rights Information about rights held in and over the resource

ircam
Centre
Pompidou

G. Peeters 44

VI. Representations of sound's description
MPEG-7

→

è MPEG-7 Multimedia Content Description Interface

- è new ISO standard [ISO/IEC 15938](#)
- è a meta-data standard, not a compression standard)
 - è Question : " How to find ? " currently: efficient text-based search, "Audiovisual data should be just as searchable as text"
 - è Descriptions: characterize, enable search/filtering, enable navigation
 - è Description generated manually or automatically
MPEG-7 : normalization of meta-data format (bit-stream), not of the extraction methods
 - è Idea: description linked to the actual A/V data (not necessarily located on the same place)
 - è Standardization: enables interoperability between metadata databases and applications

The diagram illustrates the MPEG-7 hierarchy. At the top is 'Mpeg-7', which branches into 'Description metadata' and 'Description Unit'. 'Description Unit' further divides into 'Content description' and 'Content management'. 'Content description' is split into 'Content entity' and 'Content abstraction'. 'Content entity' leads to 'Multimedia content', which includes 'Image (still region)', 'Video segment', 'Audio segment', 'Audio Visual segment', and 'Multimedia segment'. 'Content abstraction' leads to 'Summary description' and other options. Below 'Multimedia content', there are boxes for 'Media Information', 'Media locator', 'Creation Information', 'Usage Information', 'Audio descriptor', and 'Audio descriptor-scheme'. These are further detailed with 'Temporal decomposition', 'Media source decomposition', 'Matching Hint', 'Point of view', 'Relation', 'pitch spectral flatness timbre', and 'melody spoken content'. A small audio waveform and spectrogram are shown at the bottom right.

ircam
Centre
Pompidou

G. Peeters 45

VI. Representations of sound's description
MPEG-7

→

è MPEG-7 Multimedia Content Description Interface

- è Bibliographic-like information
 - è media
 - è meta (creation / production / usage)

This diagram is identical to the one above, but with red circles highlighting the 'Media Information', 'Media locator', and 'Creation Information' boxes in the lower section of the hierarchy.

ircam
Centre
Pompidou

G. Peeters 46

VI. Representations of sound's description
MPEG-7

- è MPEG-7 Multimedia Content Description Interface
 - è Bibliographic-like information
 - è media
 - è meta (creation / production / usage)
 - è Content description

ircam
Centre
Pompidou

G. Peeters 47

VI. Representations of sound's description
MPEG-7

- è MPEG-7 Multimedia Content Description Interface
 - è Bibliographic-like information
 - è media
 - è meta (creation / production / usage)
 - è Content description
 - è Structure description (table of content)

ircam
Centre
Pompidou

G. Peeters 48

VI. Representations of sound's description
MPEG-7

- è MPEG-7 Multimedia Content Description Interface
 - è Bibliographic-like information
 - è media
 - è meta (creation / production / usage)
 - è Content description
 - è Structure description (table of content)
 - è Semantic description (index)
 - è Navigation and access control
 - è summary

ircam
Centre
Pompidou

G. Peeters 49

VI. Representations of sound's description
MPEG-7

- è Part 1. MPEG-7 Systems: The tools that are needed to prepare MPEG-7 Descriptions for efficient transport and storage, and to allow synchronization between content en descriptions. Tools related to managing and protecting intellectual property
- è Part 2. MPEG-7 Description Definition Language: The language for defining new Description Schemes and perhaps eventually also for new Descriptors.
- è Part 3. MPEG-7 Visual: The Descriptors and Description Schemes dealing with (only) Visual descriptions
- è Part 4. MPEG-7 Audio: The Descriptors and Description Schemes dealing with (only) Audio descriptions
- è Part 5. MPEG-7 Multimedia Description Schemes: The Descriptors and Description Schemes dealing with generic features and multimedia descriptions
- è Part 6. MPEG-7 Reference Software: A software implementation of relevant parts of the
- è Part 7. MPEG-7 Standard MPEG-7 Conformance: Guidelines and procedures for testing conformance of MPEG-7 implementations.

ircam
Centre
Pompidou

G. Peeters 50

→ VI. Representations of sound's description
MPEG-7

è MPEG-7 Structure Description

The diagram illustrates the MPEG-7 structure description. On the left, a hierarchical tree shows an **Audio Program** containing multiple **Segment** boxes, which are further divided into **Segment** boxes, and finally into **Audio D** and **Audio DS** boxes. On the right, a detailed view of the **mpeg7:AudioSegmentType** structure is shown, including **Header**, **MediaTime**, **TemporalMask**, **AudioDescriptor**, **AudioDescriptorScheme**, **TemporalDecomposition**, **MediaSourceDecomposition**, **MediaTimePoint**, **MediaOtherTimePoint**, **MediaDuration**, and **MediaAnchor**.

ircam
Centre
Pompidou

G. Peeters 51

→ VI. Representations of sound's description
MPEG-7

è MPEG-7 Audio Low-level audio descriptors

è **Basic:**
Instantaneous waveform and power values

è **Basic Spectral:**
A log-frequency power spectrum, and spectral features including spectral centroid, spectral spread, and spectral flatness

è **Signal parameters:**
Fundamental frequency of quasi-periodic signals, and harmonicity of signals

è **Timbral Temporal:**
Log attack time and temporal centroid

è **Timbral Spectral:**
Specialized spectral features in a linear-frequency space, including a spectral centroid, and spectral features specific to the harmonic portions of signals, including harmonic spectral centroid, spectral deviation, spectral spread, and spectral variation.


è **Spectral Basis representations:**
Features used primarily for sound recognition, but generally useful as projections into a low-dimensional space to aid compactness and recognition.

è **Silence Descriptors**


ircam
Centre
Pompidou

G. Peeters 52

→ VI. Representations of sound's description
MPEG-7

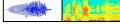


- è **MPEG-7 Audio High-level audio description tools (Ds and DSs)**
- è **Musical Instrument Timbre description tools**
Timbre descriptors aim at describing perceptual features of instrument sounds. Timbre is currently defined in the literature as the perceptual features that make two sounds having the same pitch and loudness sound different. The aim of the Timbre description tools is to describe these perceptual features with a reduced set of descriptors. The descriptors relate to notions such as "attack", "brightness" or "richness" of a sound.
- è **Sound recognition tools**
The sound recognition descriptors and description schemes are a collection of tools for indexing and categorization of general sounds, with immediate application to sound effects. Support for automatic sound identification and indexing is included as well as tools for specifying a taxonomy of sound classes and tools for specifying an ontology of sound recognizers. Such recognizers may be used to automatically index and segment sound tracks.
- è **Spoken Content description tools**
The Spoken Content description tools allow detailed description of words spoken within an audio stream. In recognition of the fact that current Automatic Speech Recognition (ASR) technologies have their limits, and that one will always encounter out-of-vocabulary utterances, the Spoken Content description tools sacrifice some compactness for robustness of search. To accomplish this, the tools represent the output and what might normally be seen as intermediate results of Automatic Speech Recognition (ASR). The tools can be used for two broad classes of retrieval scenario: indexing into and retrieval of an audio stream, and indexing of multimedia objects annotated with speech.
- è **Melody description tools (Melody, Melody contour)**
The Melody Contour DS is a compact representation for melodic information, which allows for efficient and robust melodic similarity matching, for example, in query-by-humming. The Melody Contour DS uses a 5-step contour (representing the interval difference between adjacent notes), in which intervals are quantized. The Melody Contour DS also represents basic rhythmic information by storing the number of the nearest whole beat of each note, which can dramatically increase the accuracy of matches to a query. For applications requiring greater descriptive precision or reconstruction of a given melody, the Melody DS supports an expanded descriptor set and high precision of interval encoding.




ircam
Centre
Pompidou


G. Peeters 53



→ VI. Representations of sound's description
MPEG-7



- è **MPEG-7 Audio Examples:**
 - è summary
 - è melody
 - è spectrum flatness



ircam
Centre
Pompidou

G. Peeters 54

