# Natural **transformation** of type and nature of the **voice** for extending vocal repertoire in **high-fidelity applications**

Snorre Farner,
Axel Röbel and Xavier Rodet

ircam
Centre Pompidou

Analysis/Synthesis Group
Paris, France

# Demo: Real-time voice transformation

# Overview

- ~~Real-time demonstration~~

- Introduction:
  - motivation, applications, background

- Transformation of the voice:
  - gender and age, voice quality, expressivity

- Perceptive evaluation
  of transformation of gender and age

- Conclusion

# Introduction 1

- Why transform the voice in games?
  - speech synthesis: avoid prerecorded utterances
  - voice transformation: avoid databases of many actors
  - enrichen voice repertoire for narrators and NPCs
  - design the voice of a game character based on the player's voice in multiplayer role-playing games

- Other applications:
  - educational games, e-learning, "serious games"
  - music, multimedia, audiobooks, story telling,...
  - films, dubbing, cartoon characters,...

# Introduction 2

- Ircam's objectives: artistic applications
  - music composition and composition tools

    => speech processing

    => voice and instrument transformation

- Requirements:
  - very high sound quality
  - very high degree of naturalness
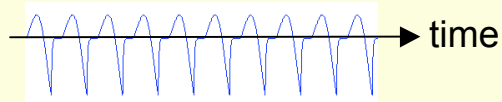  - automatic solution
  - real-time user control

ircam
Centre
Pompidou

# Voice transformation today

Two basic concepts:
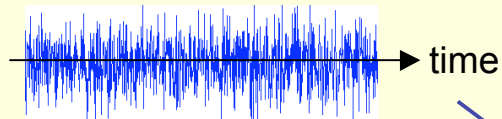
- Voice conversion: from voice A to voice B
  - often need parallel recordings of A and B
  - learning of differences between A's and B's voices
  - a new phrase of A can then be converted to B's voice
  - artifacts such as non-uniform vocal timbre

- Voice transformation: modification of general acoustic properties of the voice to transform
  - gender and age
  - voice quality: breathy voice, whispering, more or less timbred,...
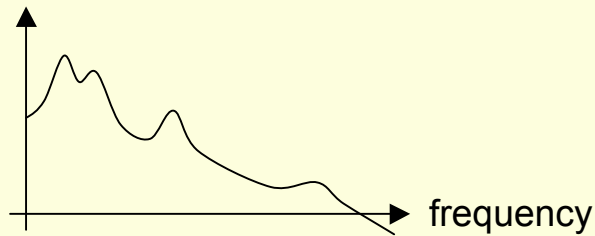  - expressivity,...

# The voice

**Pulsation of vocal folds:**

 time

**Turbulence in constrictions:**

 time

**Vocal tract resonance:**

 frequency

**Speech signal:**





vocal tract

nose cavity

mouth cavity

glottis/
vocal folds

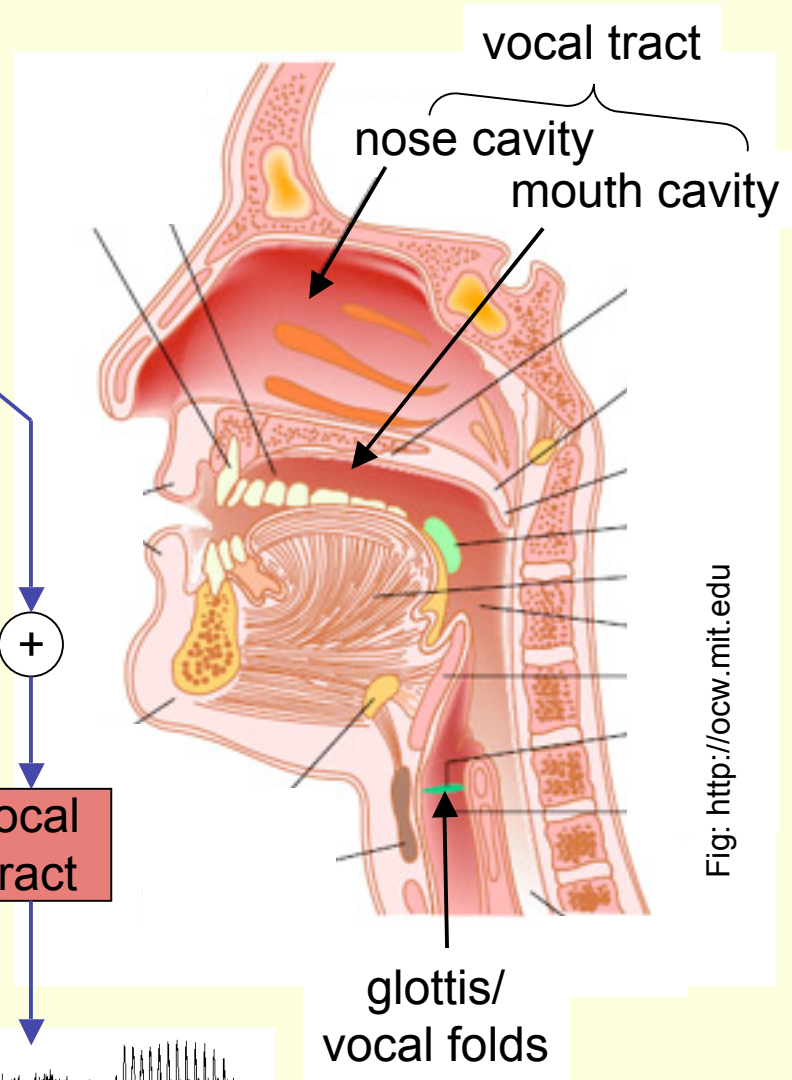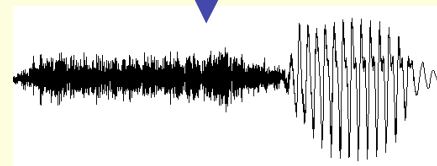Fig: http://ocw.mit.edu
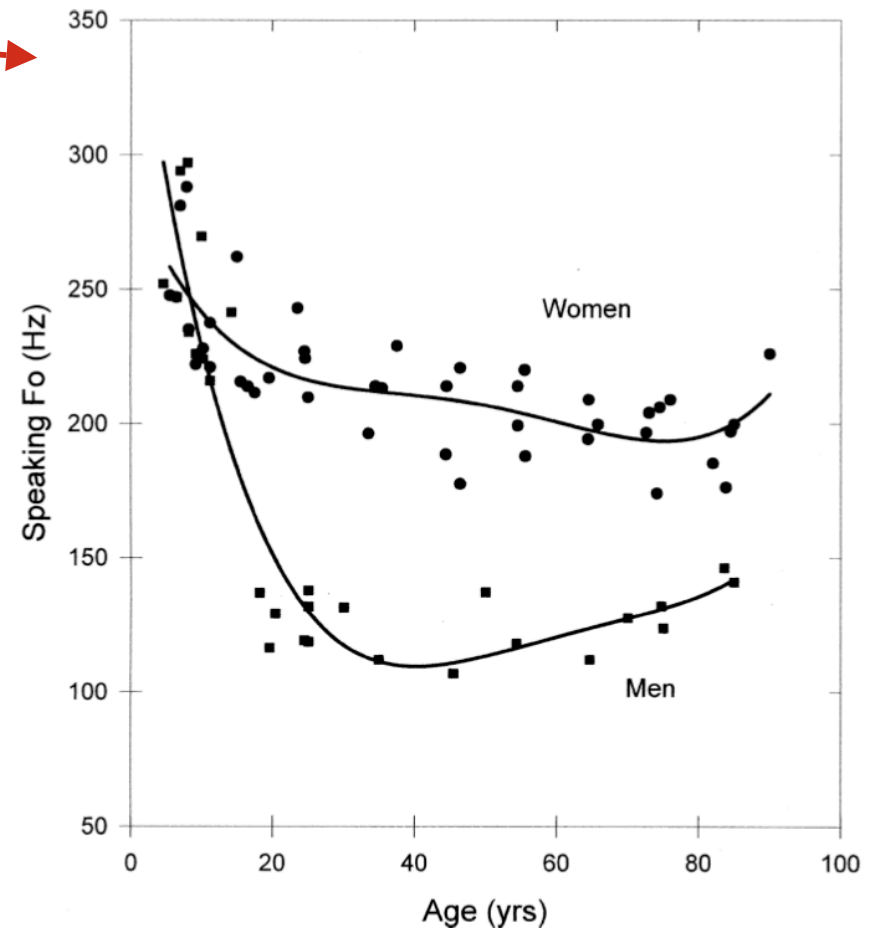
vocal tract

ircam
Centre
Pompidou

# Dependencies on gender and age

- General acoustic properties:

IrcamVoiceTrans

  - average F0
  - vocal tract; formants
  - pitch stability
  - ambitus ($F0_{max}$ - $F0_{min}$)
  - breathiness
  - speech rate

- prosody
- vocabulary
- linguistics
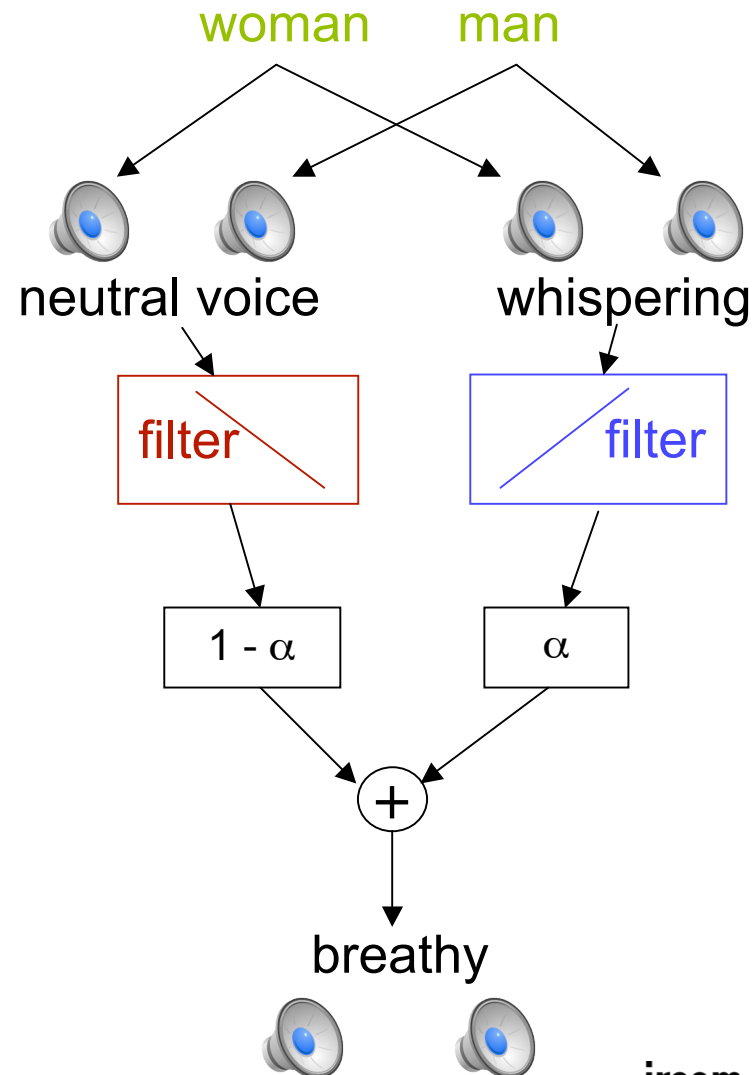


[Baken, J. Voice 2005]

# Whispering and breathy voice

## Whispering:

– filter noise by spectral envelope

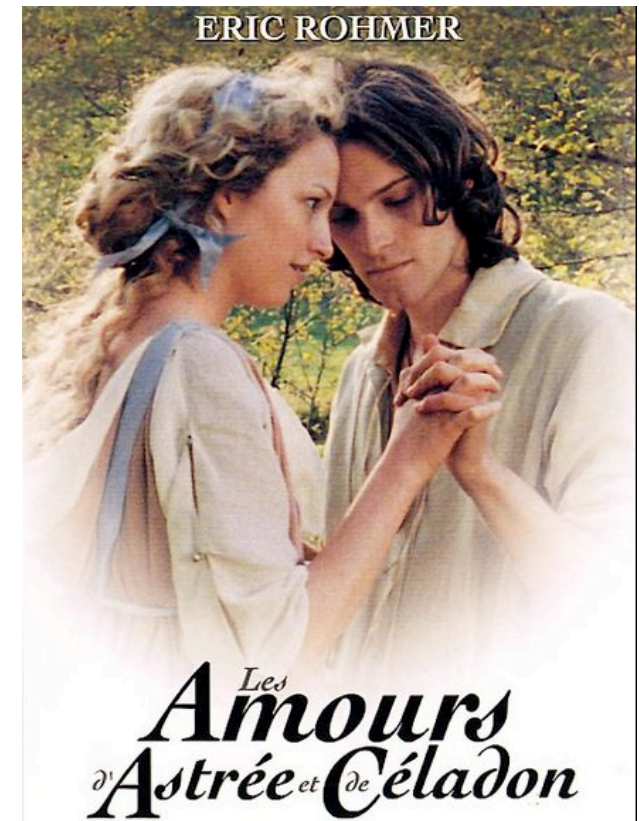– attenuate frequency bands that are voiced in original



Ex.: original: 🔊 whisper: 🔊

## Breathy voice:

# Transformation of identity

- Disguising man to woman:
  - ...also the voice: 🔊 → 🔊
  - Céladon 🔊 → Alexie 🔊

- Monologue → dialog
  🔊            🔊

- One actor to 12 persons:
  - 🔊 → 🔊 5th Blind (woman)
  - 🔊 → 🔊 Oldest Blind Woman
  - 🔊 → 🔊 Oldest Blind Man
  - 🔊 → 🔊 3rd Blind (man)



ERIC ROHMER

Les Amours d'Astrée et de Céladon



« Deux Songes
de Maeterlinck d'après Brueghel »
by J. B. Barrière, 2007

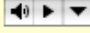# Examples: transformation of voice quality

- original: 🔊

- breathy: 🔊
- whispering: 🔊
- creaky: 🔊

- trembling: 🔊
- pitch ambitus: greater: 🔊   smaller: 🔊   zero (robot): 🔊
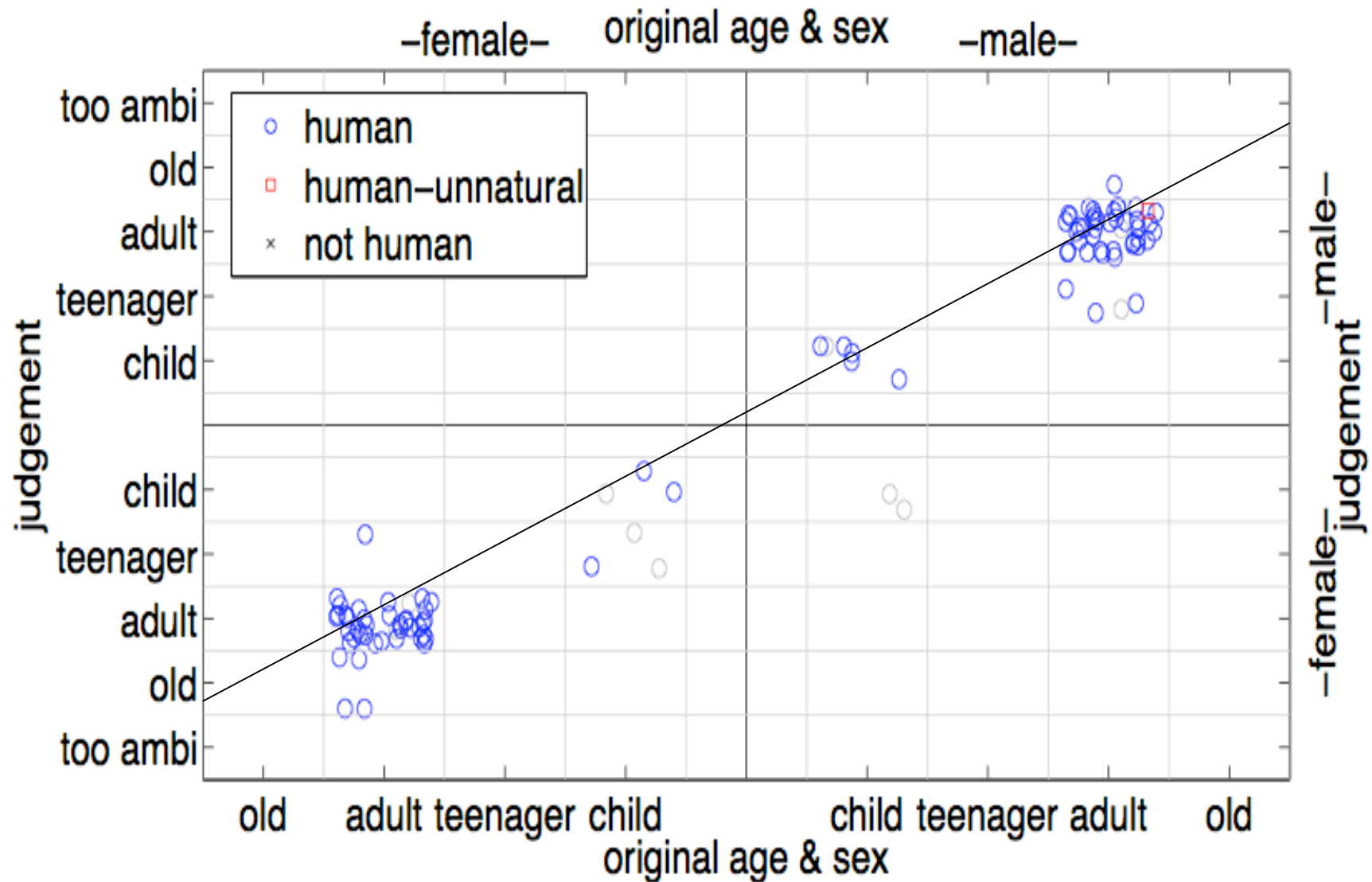- speech rate: faster: 🔊   slower: 🔊   faster vowels: 🔊

combinations
- dull: 🔊        excited: 🔊
- drunk: 🔊

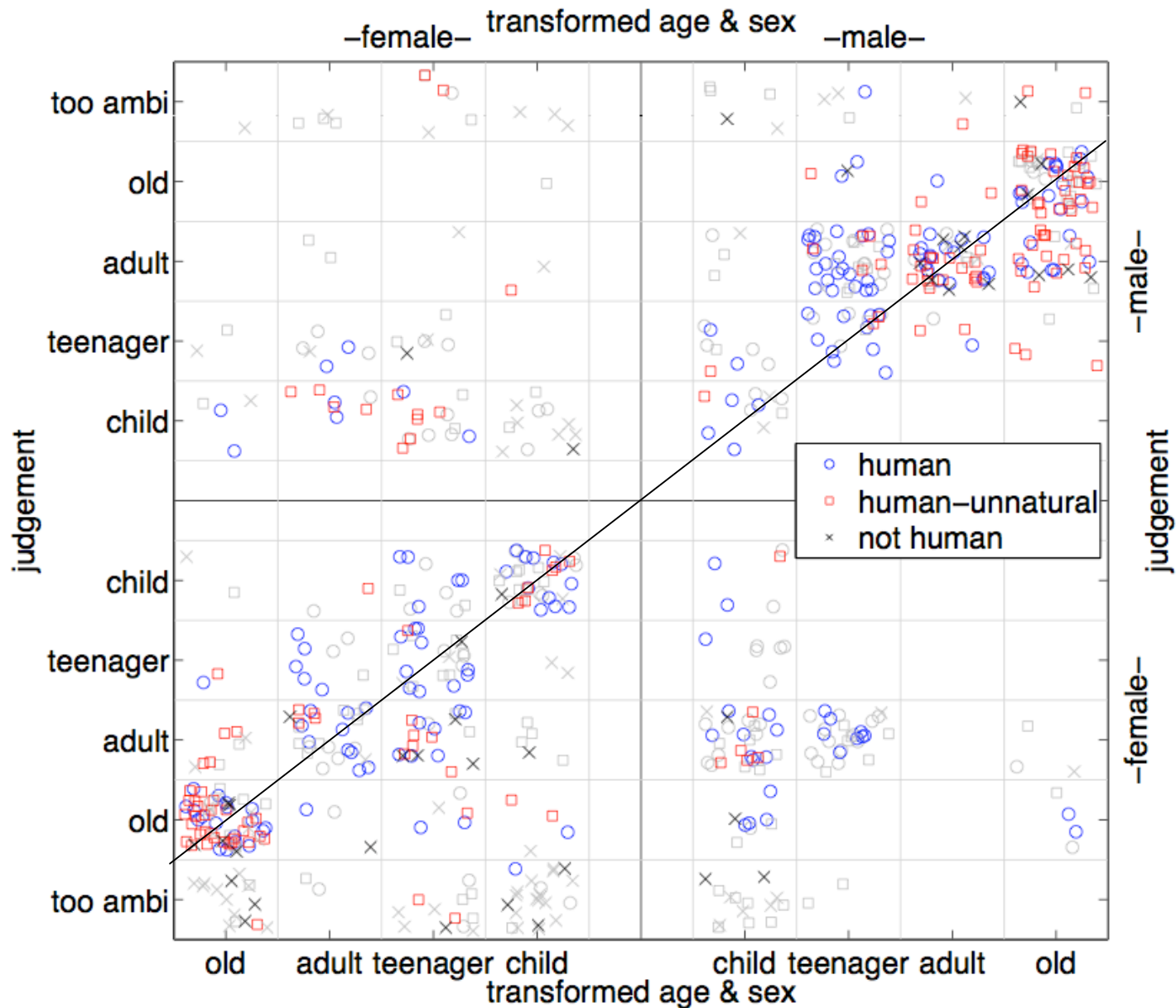ircam
Centre
Pompidou

# Perceptual evaluation

- **13 voices** (5 women, 6 men, 1 girl and 1 boy)
- **2 sentences** of 2 to 3 seconds
- **7 transformations** (male/female x 4 ages) +**original**
- **31 subjects** listening to each sentence once (26)

| The voice | | | The sound quality |
|---|---|---|---|
| Listen: ◀ ▶ ▼   Remaining samples: | | | Listen again: ◀ ▶ ▼ |
| What's the **sex** of the voice? | What's the **age** of the voice? | Does the **voice** sound like a human? | Did you notice any artefacts (buzz, echo, strange sounds/noises, etc.)? |
| ○ male<br>○ probably male<br>○ probably female<br>○ female | ○ child<br>○ teenager<br>○ adult<br>○ old<br>○ too ambiguous to tell<br><br>If uncertain, try to pick the closest. | ○ Yes, a human speaking naturally<br>○ Yes, a human speaking in an unnatural way<br>○ No, not human | ○ No<br>○ Yes, but not annoying<br>○ Yes, slightly annoying<br>○ Yes, annoying<br>○ Yes, very annoying |
| | | | NEXT |

# Evaluation: original voices

# Evaluation: transformed voices

# Examples of
# transformation of gender and age

source voice

| target voice | woman | man |
|---|---|---|
| original | . 🔊 | . 🔊 |
| little girl | . 🔊 | . 🔊 |
| teenage girl | . 🔊 | 🔊 |
| woman | | . 🔊 |
| aged woman | 🔊 | 🔊 |
| boy | 🔊 | 🔊 |
| teenage boy | . 🔊 | . 🔊 |
| man | . 🔊 | |
| aged man | . 🔊 | 🔊 |

**ircam**
Centre
Pompidou

# Conclusion

- **Perceptual evaluation:**
    - listening test focuses on artifacts
    - in real world, attention is distracted by background sounds, visual input, story line, etc.
    - transformation of pitch and timbre not enough, e.g., girl 🔊 → man 🔊 :  a girl's way of speaking

- Voice transformation is already used in high-fidelity applications: music, film, theatre

- IrcamVoiceForger:
    - C++ library
    - real-time

# Demo: one actor → 4 characters



Produced with
Living Actor™ character
and voices by IRCAM

[Characters and animation by Cantoche]

Original: 🔊