

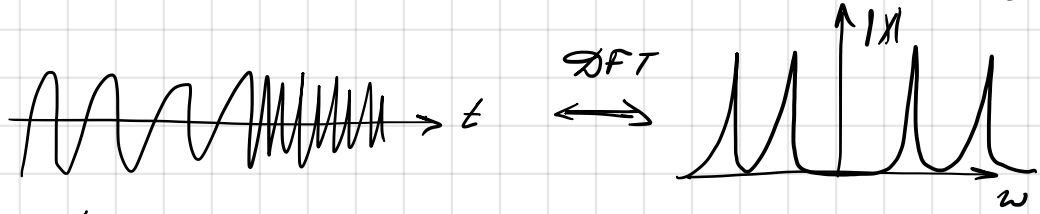
Современные методы  
распознавания и синтеза речи  
Лекция 3

# Лекция 3 Частотно-временной анализ

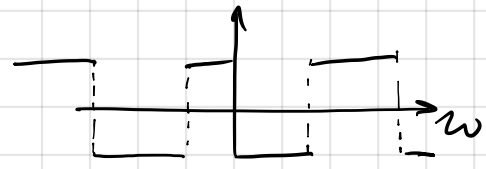
## §1 Short-time Fourier transform

Преобразование Фурье показывает "что" находится в сигнале, но не "где". Например, DFT of музыки  $\Rightarrow$  узнаем все ноты, коэф-ты перемись, но конкретной моменту знать не сумеем.

Пример



Временная область  
(Time domain)  
[сек.]



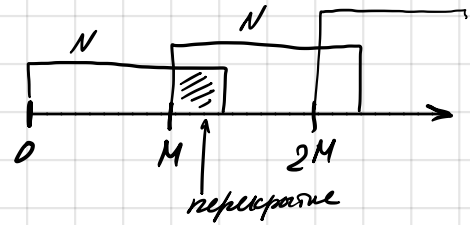
Частотная область  
(Frequency domain)  
[1/сек.]

Спектрограмма:

$$S[k, m] = \sum_{i=0}^{N-1} x[mM+i] e^{j \frac{2\pi}{N} \cdot ik}$$

↑            ↑  
частота    время

$$S[k, m] = W_N \begin{bmatrix} x[0] & x[M] & \dots \\ x[1] & x[M+1] & \dots \\ x[2] & x[M+2] & \dots \\ \vdots & \vdots & \ddots \\ x[N-1] & x[N+M-1] & \dots \end{bmatrix}$$

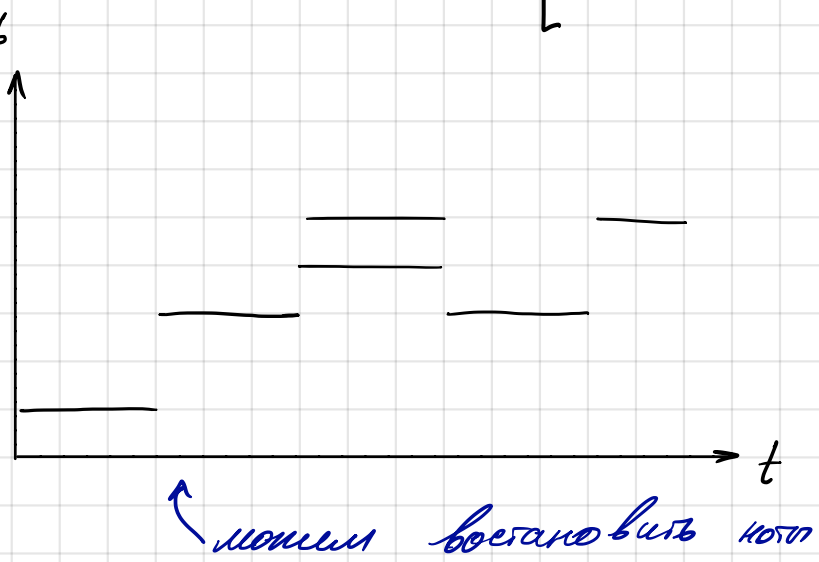


Принцип неопределенности:

не можем одновременно точно узнать частоту и время.

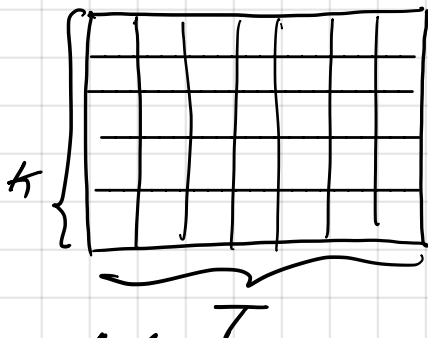
$N \uparrow \Rightarrow$  лучше разрешение по частотам, но больше время

$N \downarrow \Rightarrow$  хуже частотное разрешение, но лучше локализация по времени



Применение: выделение признаков из звука.

- обходясь на исходном сигнале фрунго: небольшая задержка или изменение громкости  $\Rightarrow$  сильное отличие сигнала
- спектрограмма более робастна к таким изменениям



ML  $\rightarrow$

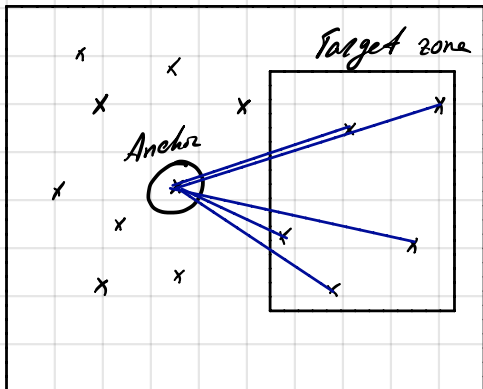
- детекция ключевых фраз
- идентификация голоса
- поиск кода фразы

## Audio fingerprinting (Shazam)

1. Выделяем локальные пики в спектрограмме. Эти признаки робастны к сильному шуму.

- > Фильтр высоких частот
- > Зад. max в окрестности

2. Выделение отрезков (fingerprinting)



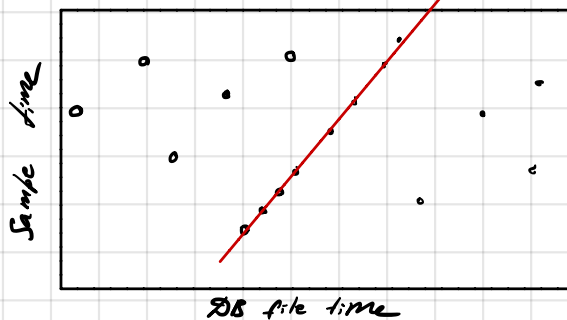
Для каждой пары:  $[f_{\text{anc}}, f_{\text{tg}}, \Delta t]$

Снижение памяти: квантование  $f_{\text{anc}}, \Delta t$ , запись одним числом, SHA где повышение энтропии (более точное)

+ сократим  $f_{\text{anc}}$

совпадение!

3. Поиск аудиофрагмента:

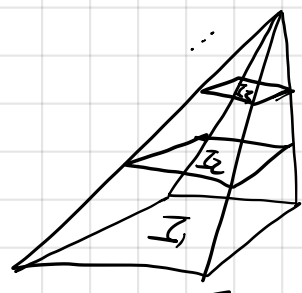
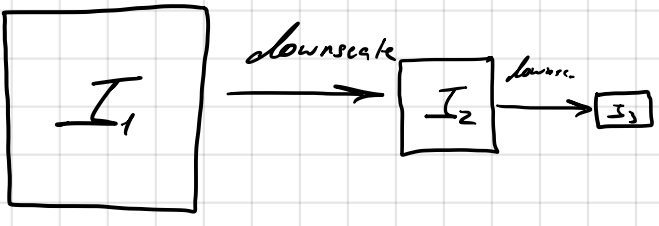


## §2 Multi-resolution & subband coding

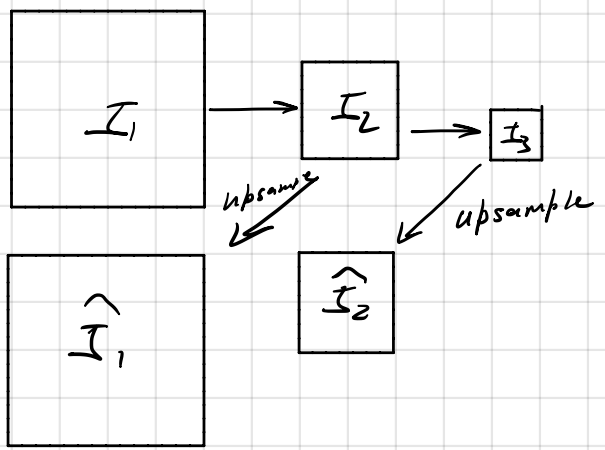
Идея: выделим широкополосные компоненты, удалим их из сигнала. Выделим более высокие частоты, удалим их, etc.

Последенно расчлениваем сигнал, разделив его на компоненты.

Пример Хаусова пирамида

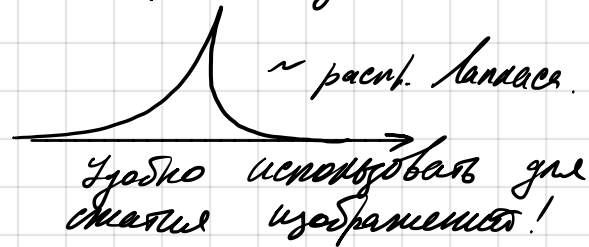


Получили представление на разных уровнях. Посмотрим на ошибку реконструкции

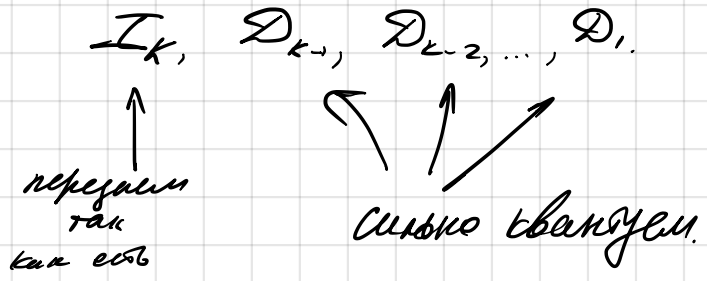


$D_k = I_k - \text{upsample}(I_{k+1})$

ошибка реконструкции



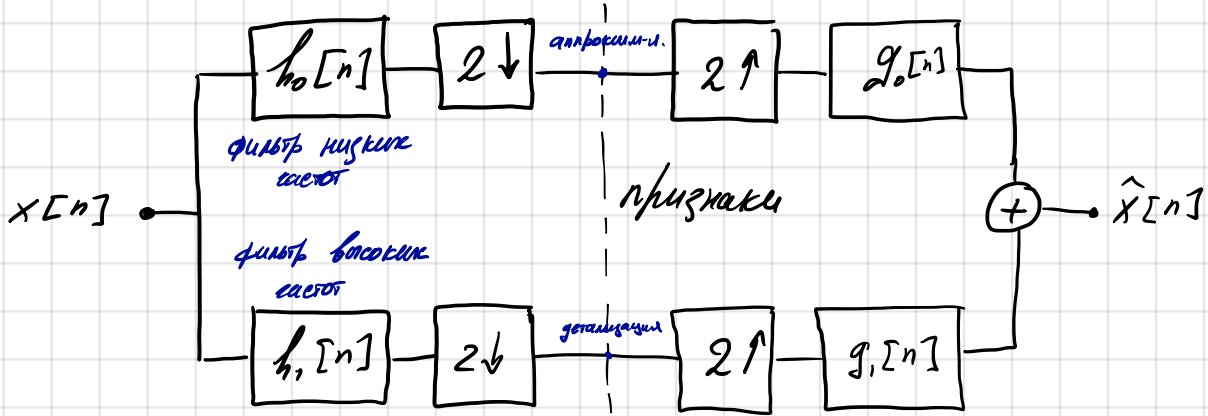
Пирамида Лангаса:



Техника 'multiresolution analysis'.

Subband coding

Разделим на низкочастотную (аппроксимация) и высокочастотную (детализация) компоненты.



Как получить две версии  $h_0, h_1, g_0, g_1$ ? Хотим избежать ошибок.

$$X(z) = \sum_{n=-\infty}^{+\infty} x[n] z^{-n}$$

$$x_{\text{down}}[n] = x[2n] \iff X_{\text{down}}(z) = \frac{1}{2} (X(z^{1/2}) + X(-z^{1/2}))$$

$$X_{\text{down}}(z) = \sum_{n=-\infty}^{+\infty} x[2n] z^{-n} = \sum_{n=-\infty}^{+\infty} x[2n] z^{-n} = \sum_{n=-\infty}^{+\infty} x[2n] (z^{1/2})^{-2n} = \frac{1}{2} (X(z^{1/2}) - X(-z^{1/2}))$$

$$x_{\text{up}}[n] = \begin{cases} x[n/2], & n=0,2,4,\dots \\ 0, & \text{otherwise} \end{cases} \iff X_{\text{up}}(z) = X(z^2)$$

$$\hat{X}[2] = \frac{1}{2} (X(z) + X(-z))$$

$\underbrace{\hspace{100px}}$   
исходный сигнал
 $\underbrace{\hspace{100px}}$   
искаженная версия

Хотим избежать искажений

$$\hat{X}(z) = \frac{1}{2} G_0(z) [H_0(z)X(z) + H_0(-z)X(-z)] + \frac{1}{2} G_1(z) [H_1(z)X(z) + H_1(-z)X(-z)] =$$

$$= \frac{1}{2} \underbrace{[G_0(z)H_0(z) + G_1(z)H_1(z)]}_{=2} X(z) + \frac{1}{2} \underbrace{[G_0(z)H_0(-z) + G_1(z)H_1(-z)]}_{=0} X(-z)$$

$$[G_0(z) \quad G_1(z)] \begin{bmatrix} H_0(z) & H_0(-z) \\ H_1(z) & H_1(-z) \end{bmatrix} = [2 \quad 0]$$

$$H_m \quad H_m^{-1} = \frac{1}{\det H_m} \begin{bmatrix} H_1(-z) & -H_1(z) \\ -H_0(-z) & H_0(z) \end{bmatrix}$$

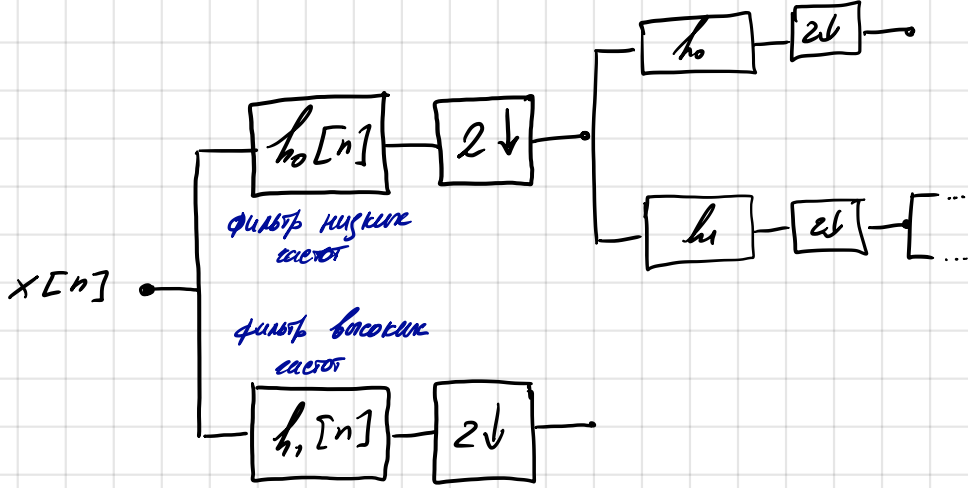
$$\begin{bmatrix} G_0(z) \\ G_1(z) \end{bmatrix} = \frac{2}{\det(H_m(z))} \begin{bmatrix} H_1(-z) \\ -H_0(-z) \end{bmatrix}$$

где FIR
 $d \cdot z^{-(2L-1)}$

Интересно что  $H_1$  орт-г  $G_0$ ,  
 $H_0 - G_1$

Высокочастотную компоненту можно пожелать реализовать: чтобы разделить на несколько полос частот.

Применение: разное кодирование для разных полос



### §3 Вейвлет преобразование

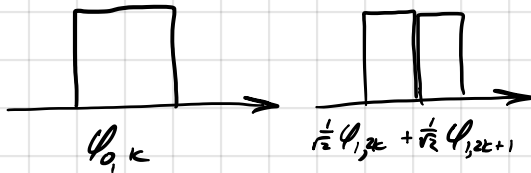
Пусть  $\varphi(x)$  — квадратично суммируемая функция. Применим к ней целочисленные сдвиги и растяжения:

$$\varphi_{J,k}(x) = 2^{J/2} \varphi(2^J x - k)$$

$\varphi(x)$  — scaling function / father wavelet.

Порядковые пространства:  $V_J = \text{Span}_k \{ \varphi_{J,k}(x) \} \subset L_2(\mathbb{R})$

Пример:  $\varphi(x) = \begin{cases} 1, & x \in [0, 1) \\ 0, & \text{иначе} \end{cases}$   
 scaling function Хаара

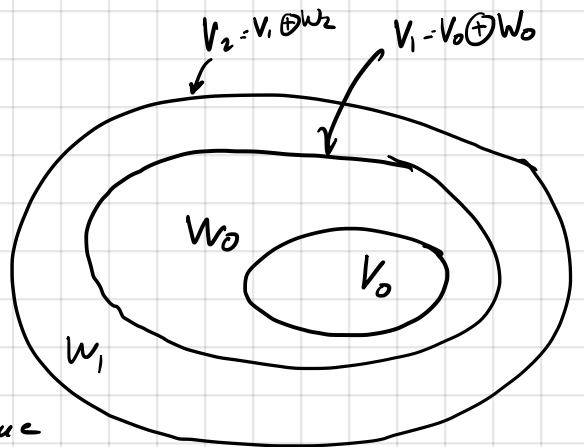


$$\varphi_{0,k}(x) = \frac{1}{\sqrt{2}} \varphi_{1,2k}(x) + \frac{1}{\sqrt{2}} \varphi_{1,2k+1}(x)$$

$$V_1 \subset V_2 \subset V_3 \subset \dots \subset L_2(\mathbb{R}) = V_\infty$$

Будем раскладывать функции по парам:

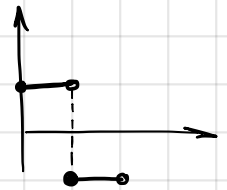
$$L_2(\mathbb{R}) = V_0 \oplus W_1 \oplus W_2 \oplus \dots$$



Как построить базис  $W_k$ ? Мы знаем базис  $V_k, V_{k+1}, V_{k+1} = V_k \oplus W_k$ . Ортогонализация Грамма-Шмидта!

$$\langle \varphi_{J,k}(x), \varphi_{J,l}(x) \rangle = 0 \quad \forall J, k, l \in \mathbb{Z}$$

Для scaling-функции Хаара:  $\psi(x) = \begin{cases} 1, & 0.5 \leq x < 1 \\ -1, & 0 \leq x < 0.5 \\ 0, & \text{иначе} \end{cases}$



$$\psi_{J,c}(x) = 2^{J/2} \psi(2^J x - c)$$

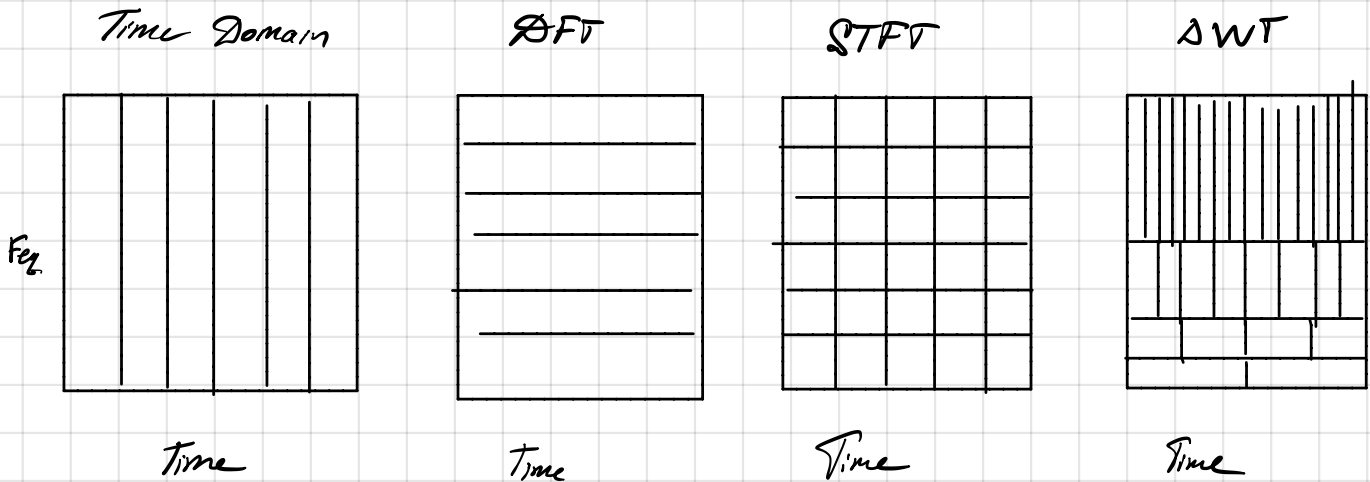
материнский вейвлет

Разложение:  $f(x) = \underbrace{\sum_k c_{J_0}(k) \psi_{J_0, k}(x)}_{\text{аппроксимация}} + \underbrace{\sum_{J=J_0}^{+\infty} \sum_k d_J(k) \psi_{J, k}(x)}_{\text{детализация}}$

$$c_{J_0} = \langle \psi_{J_0, k}(x), f(x) \rangle$$

$$d_J(k) = \langle \psi_{J, k}(x), f(x) \rangle$$

По аналогии с DTFT: кент. вейвлет преобразование:  $\psi_{s, \tau}(x) = \frac{1}{\sqrt{s}} \psi\left(\frac{x-\tau}{s}\right)$



Применение: стазис (JPEG-2000), шумоподавление.

Анализ: CWT показывает сходство участка сигнала с вейвлетом и промасштабированному вейвлету.