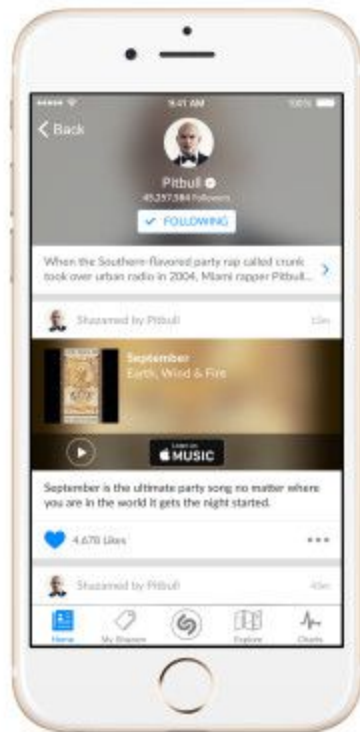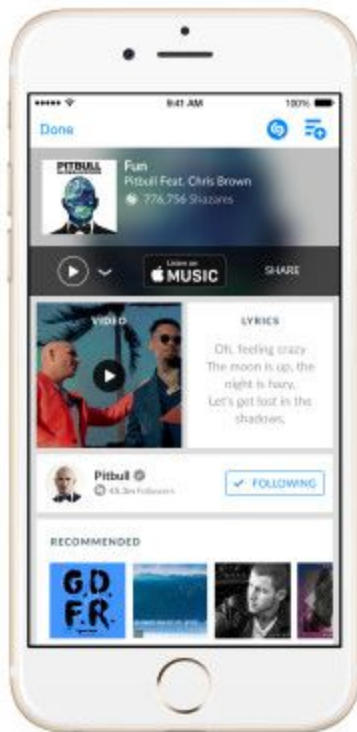# An industrial-strength- audio search algorithm

Hello, I'm Tom

# What is Shazam?

- Identifies **exact tracks** of music.
- Only needs small samples (seconds)
- Robust to noise

# What isn't Shazam

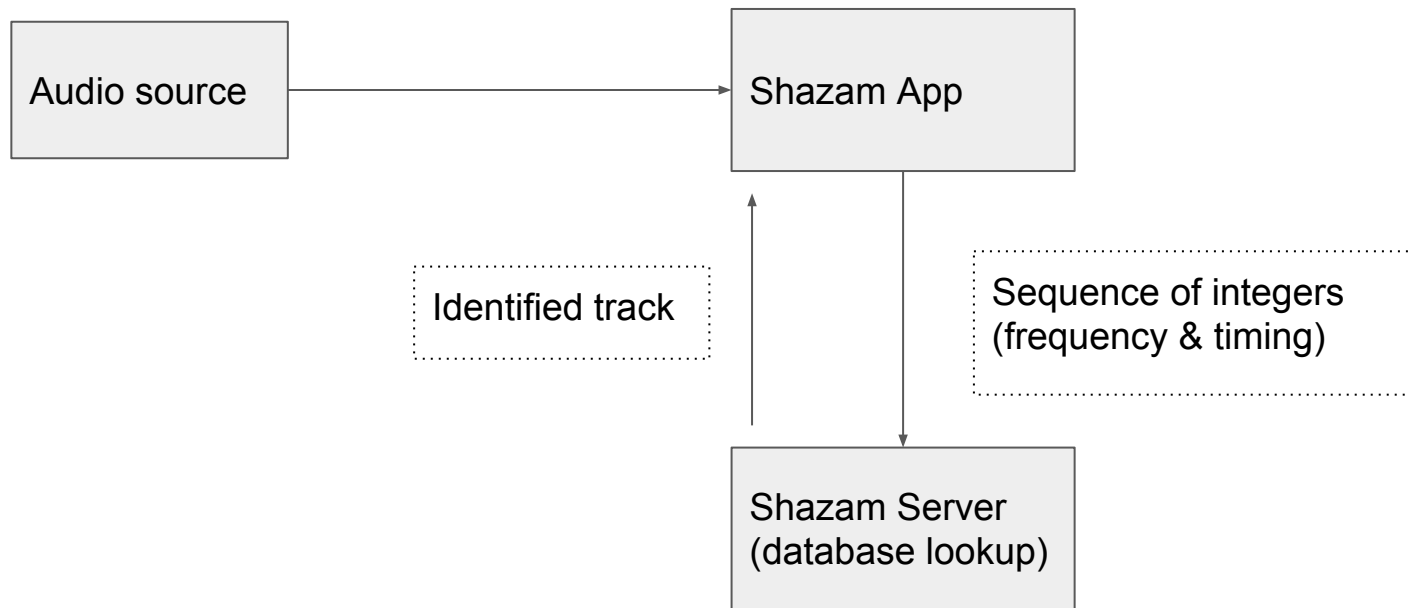- Not designed to detect live recordings.

# Shazam is **old**

Introduced in 1999!

# Basic idea: audio fingerprinting

Two key pieces to Shazam:

1) Construction of "fingerprints"
   a) Contain frequency and timing information
2) Lookup of fingerprints

# Part 1: Construction of the fingerprints
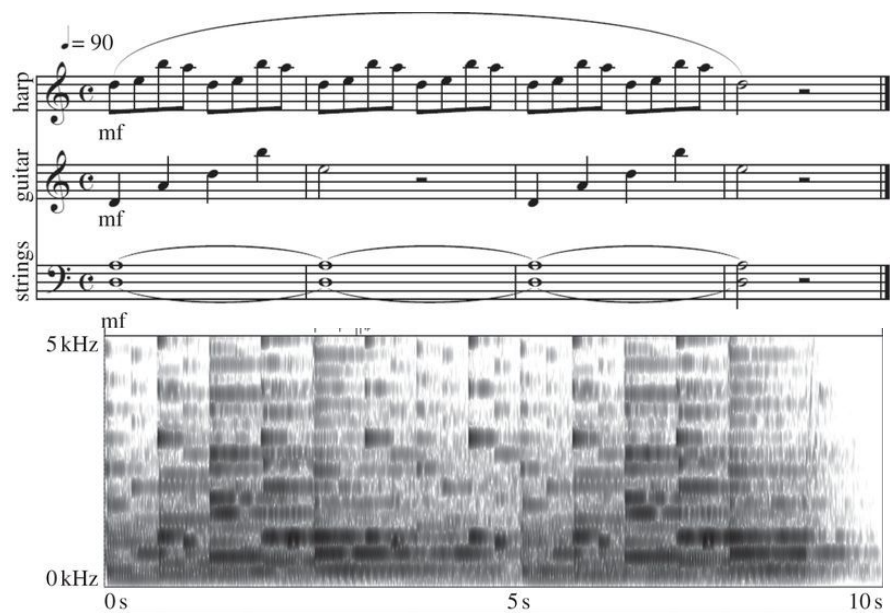
# Spectrograms and Sheet music

# Spectrograms and Sheet music

# Piano C4 note (~260Hz fundamental freq)

# How to select interesting frequencies?



Spectrogram of "Blurred Lines" by Robin Thicke

# Peak detection



Spectrogram of "Blurred Lines" by Robin Thicke

- Extremely robust to noise
- Highly reproducible

# Open source implementation!

https://github.com/worldveil/dejavu

Uses scipy's `maximum_filter` for peak detection.

# Constructing the hashes 1: quantization

Frequencies are binned into 1024 values


=> We only need 10 bits to encode a quantized frequency.

# Constructing the hashes 2: a wrong idea

What if we sent off the locations of the peaks?

In other words, send (quantized)

$$\texttt{(time\_offset, frequency)}$$

**What's wrong with this?**

# Lookup on the database is the problem

- We can't key off the pair `(time_offset, frequency)`: database would be enormous, and processing would be terrible.
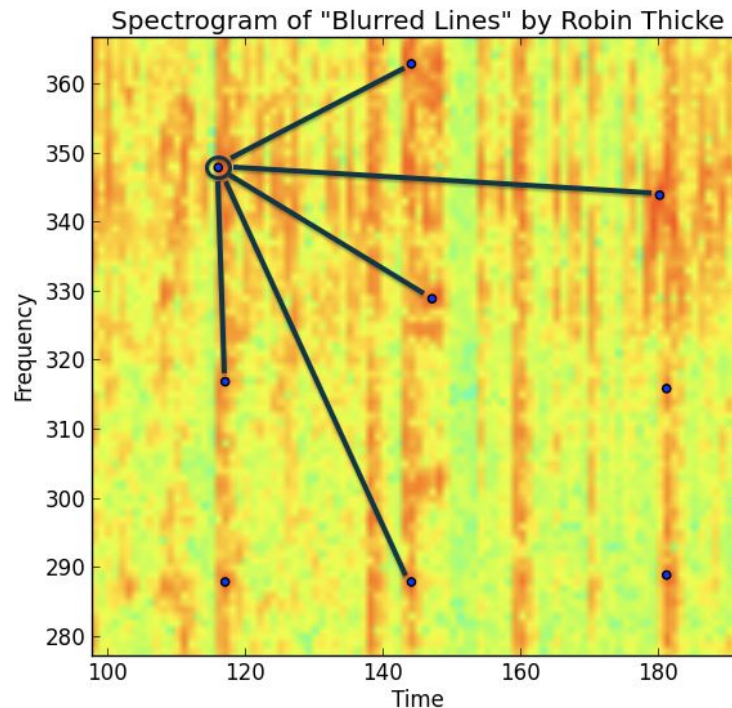- Frequency alone leads to many prospective matches.

# Shazam's solution: look at frequency pairs

Anchor: `(t0, f0)`
Target: `(t1, f1)`

Hash is 32-bit integer of:

`[10 bits f0,`
`10 bits f1,`
`10 bits (t1 - t0)]`



Spectrogram of "Blurred Lines" by Robin Thicke

# Server side:

Only need a linear scan of each track to generate fingerprints

# Part II: Looking up fingerprints

# Server side: lookup

Incoming stream: h0:t0, h1:t1, h2:t2, …

  (recall each hash = [freq0, freq1, time_delta])

Form buckets:

Song_xyz: h1:t1_xyz, h4:t4_xyz. h7:t7_xyz

Song_abc: h0:t0_abc, h1:t1_abc, h3:t3_abc, h5:t5_abc, h6:t6_abc

Song_123: h0:t0_123, h8:t8_123

# Server side: computing correlations
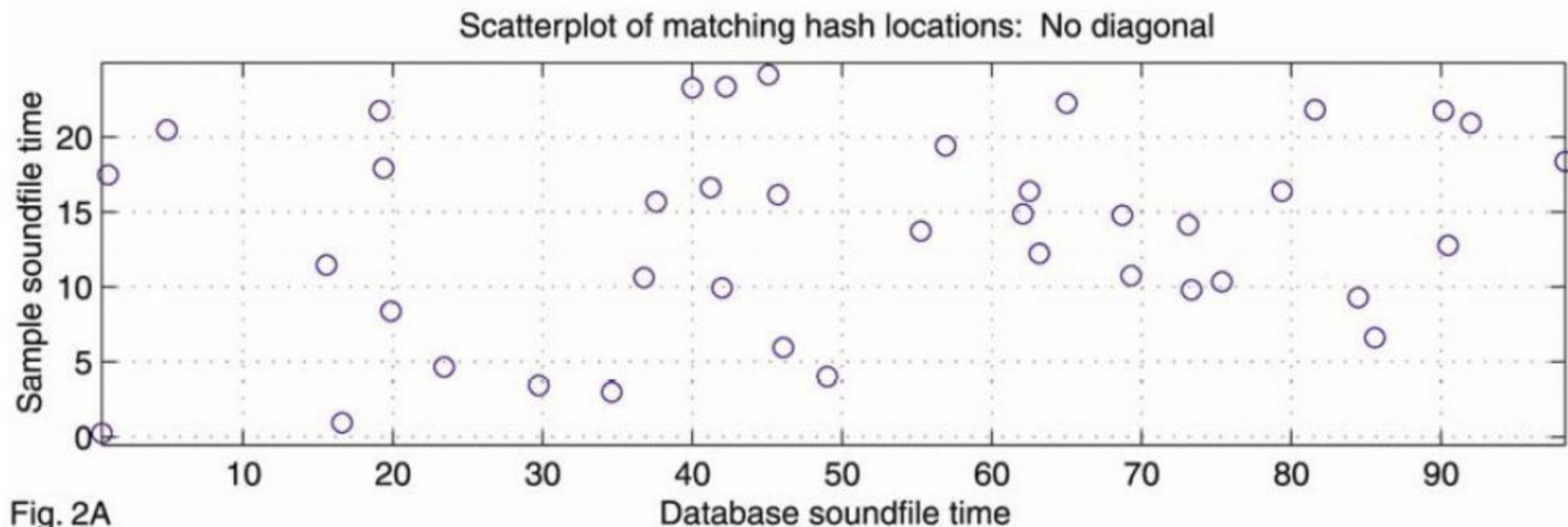
Bad match: many key matches but not time-aligned



Scatterplot of matching hash locations: No diagonal

Fig. 2A

# Server side: computing correlations

Good match: many are time-aligned



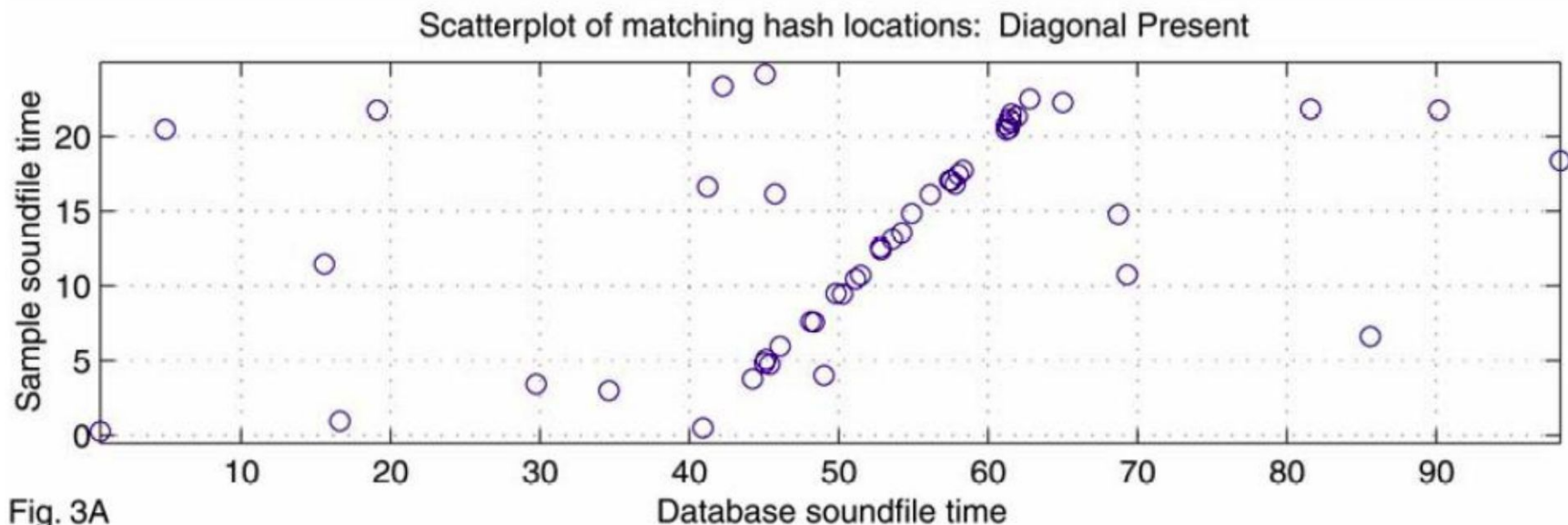Scatterplot of matching hash locations: Diagonal Present

Fig. 3A

# How does shazam measure correlations?

- Could use robust regression, R^2 or whatnot (time complexity anyone?)
- Much simpler approach: histograms (time complexity anyone?):
  - Denote $\{\ t_i\ \}$ set of time offsets from sample, $\{\ t'_i\ \}$ time offsets from database.
  - If from same song, $t_i\ =\ t'_i\ +\ c$ for some constant $c$.
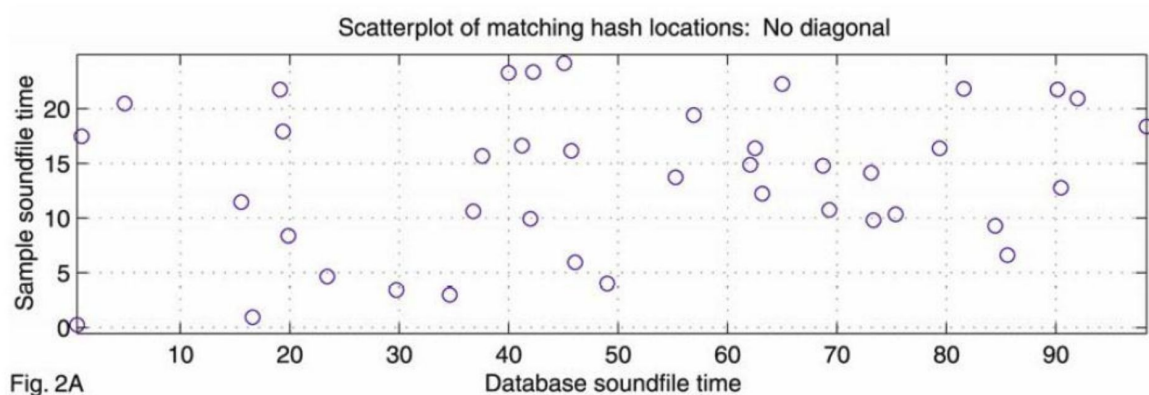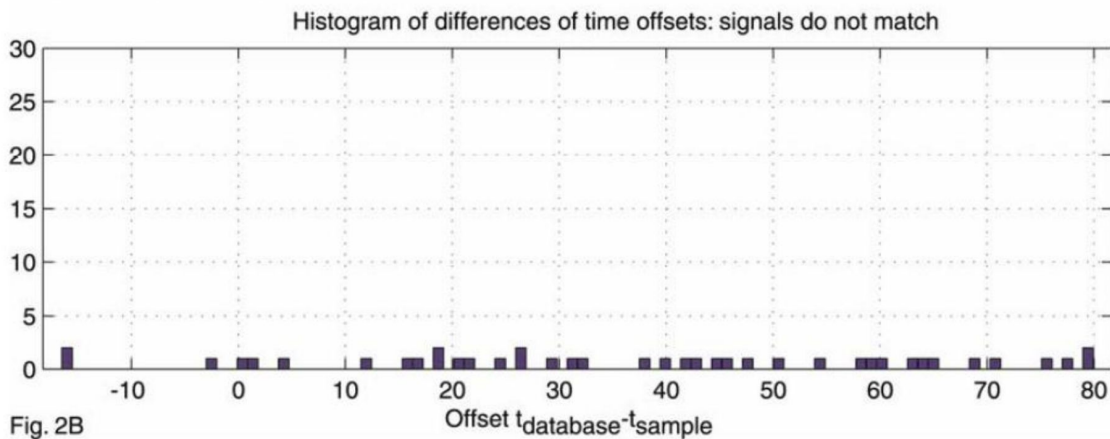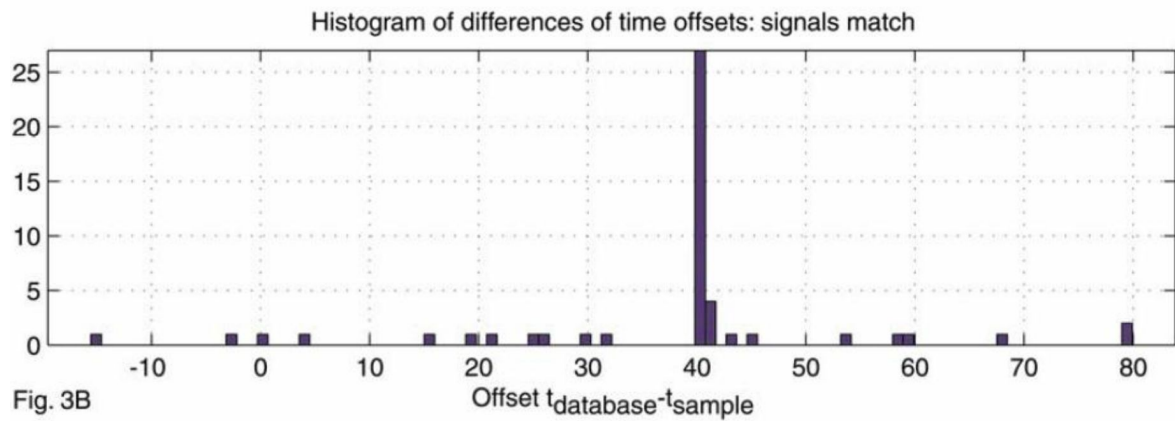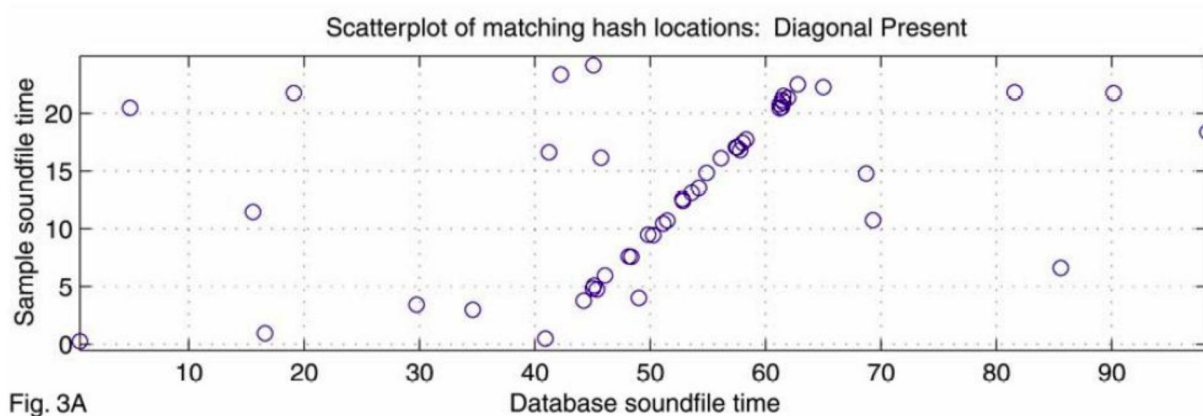  - Form histograms of $\{\ t_i\ -\ t'_i\ \}$ and look for peaks.

Fig. 2A

Fig. 2B

Scatterplot of matching hash locations: Diagonal Present

Fig. 3A

Histogram of differences of time offsets: signals match

Fig. 3B

# Questions? Thank you!

[thomas.d.peters@gmail.com](mailto:thomas.d.peters@gmail.com)

- "An industrial-strength audio search algorithm", by Avery Li-chun Wang. *Proceedings of the 4 th International Conference on Music Information Retrieval*
- [https://github.com/worldveil/dejavu](https://github.com/worldveil/dejavu)