# corti

Revolutionizing healthcare.
One patient consultation at a time.

**Lars Maaløe**
Co-Founder & Chief Technology Officer @ Corti
Adj. Associate Professor @ Technical University of Denmark

corti

**Digital health is exploding**

There's now 50 billion patient consultations conducted a year.

corti

**Only 1% of patient consultations are quality assured**

The pressure on caregivers is so high that **88% of diagnoses are altered** if tested by a second opinion.
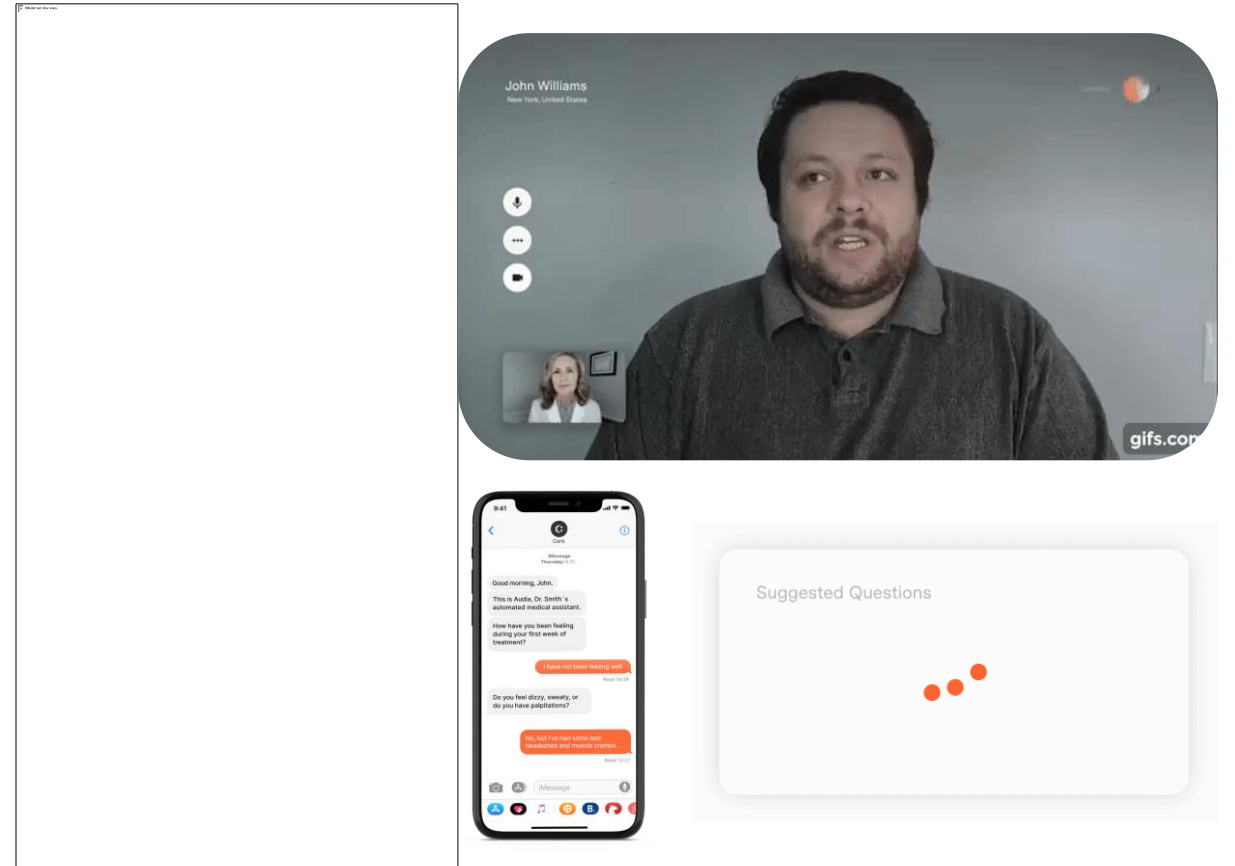
corti

**Physicians are becoming data clerks**

# Telehealth physicians spend up to **50% of their day** in EHR documenting consultations.

**Intelligent augmentations**

**Corti** has built an artificial intelligence that **listens** in and **understands** virtual consultations.



corti

# Our AI is built by a world leading team of engineers in close collaboration with academia

**+369** h-index   **+938** i-index   **+248k** citations

Accumulated over 6x Professors, 2x Associate Professors, and a
number of PostDocs and PhD scholars

*Excerpt of some of our latest work*



*In direct competition with the major AI labs*

Prior to this presentation, our AI has learned from

# +15,000,000

patient calls to medical command centers.

corti

**2016, we started where urgency is highest**

Emergencies are our beachhead market, helping where it matters the most and getting the best data.

**Servicing +50 million yearly encounters**

We service customers in Europe and USA with some of the best annotated healthcare data, <span style="color:orange">continuously improving the AI.</span>

corti

# Our technology can solve the biggest challenges for healthcare services



## Healthcare worker

*Burnout and Churn*



## Citizens

*Healthcare Quality*



## Organization

*Trust and Credibility*

1 *"EMS services warn of 'crippling labor shortage' undermining 911 system."* NBC News. 8 Oct. 2021, https://nbcnews.to/3M2PtAp
2 *"Stress on the front lines of covid-19 - The Washington Post."* 6 Apr. 2021, https://wapo.st/3GRjnUy
3 *"The longitudinal study of turnover and the cost of turnover in EMS."* 11 June 2010, https://bit.ly/3v8DiMz

corti

# Corti analyzes **100%** of your communication



**Automatic annotations of all communication**

What we are most commonly known for

# Triage a patient through workflow software with intelligent decision support



*Link to video*

corti

# Reduced call duration in Sweden by +20%



"Our more than 800 operators use Corti for safer and faster medical triaging. This has allowed us to significantly increase patient safety through protocol adherence, while reducing average call duration by more than 20 percent and counting."

Jannice Mattsson, COO, SOS Alarm

*_Link to video_

corti

Following the patient journey

# Since 2016, we have optimized towards supporting the entire healthcare value chain

corti

Jill, a Seattle resident, feels terrible. **She calls 911.**

Joe, the dispatcher, does not consider this an acute emergency. **He redirects to a nurse line.**

Jane, the Nurse in Dallas, finds it necessary to arrange for a Telemedicine call.
**She arranges a MD call.**

Brian, the Nurse in Inglewood, finds it necessary for a physical consultation. **She arranges the consultation.**

Jannet, the Doctor in Seattle, finds the need for a procedure.
**She arranges the procedure.**

Donald, the Surgeon in Seattle, performs the procedure.
**He wishes Jill all the best.**

corti

Example: Alleviating the administrative burden

# Corti Code documents the patient interaction



corti

# How it works



corti

Recent Example: Winning a +100 hospital health provider

# ICD-10 coding from audio in competition with a large set of companies



Corti Code
AI-Powered Medical Coding

Executive and Technical Report

Davide Paganini · Lead Product Manager. Alexander Junge · Data Science Lead.
Lorenzo Belgrano · Machine Learning Engineer. Sotiris Lamprinidis · Machine Learning Engineer.
Henrik Cullen · VP of Product. Lars Maaløe · Chief Technology Officer.

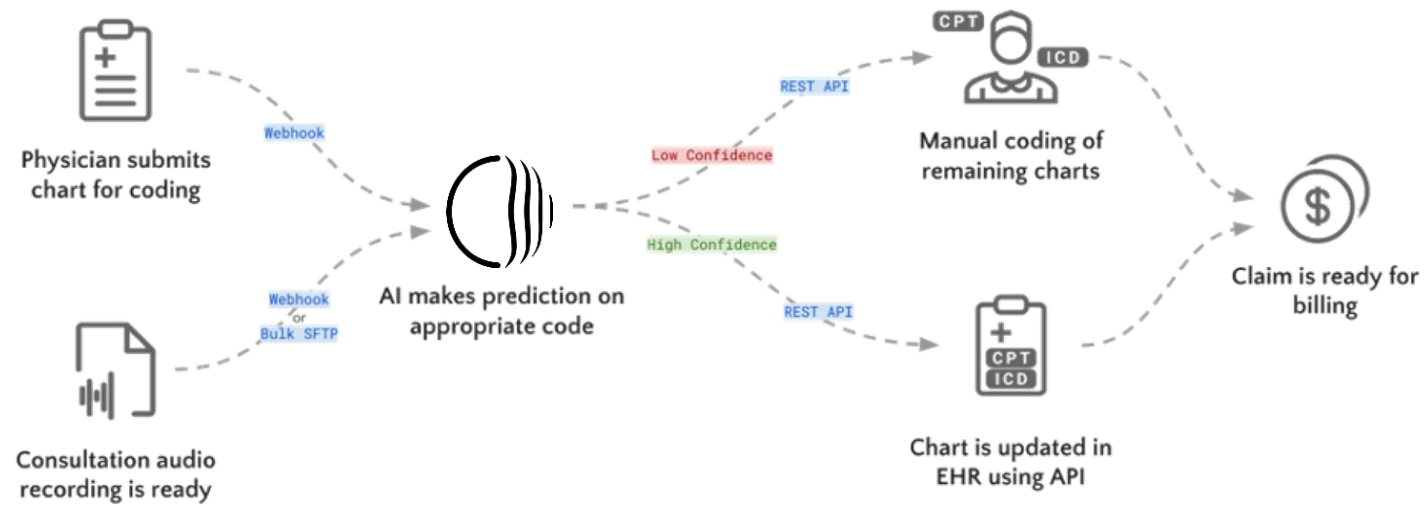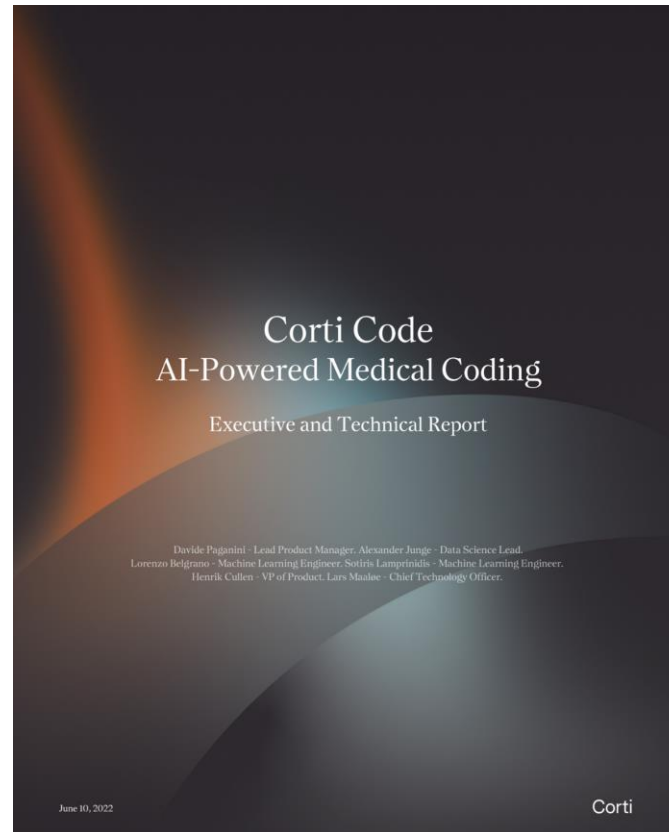June 10, 2022

Corti

-> +75% of consultations are fully automated by AI.

-> +95% of consultations are gets the right code(s) through a top-5 recommendation engine.

corti

Contact
Lars Maaløe
Co-founder & CTO
lm@corti.ai