

Provenance-based Data Analytics CI for High-frequency Mobile Sensor Data (mProv)

Emre Ertin (Ohio State), Zach Ives (UPenn), Santosh Kumar (Memphis),
Jim Rehg (Georgia Tech), Ida Sim (UCSF), Mani Srivastava (UCLA)



*NIH Big Data to
Knowledge (BD2K)*

Growing Potential of High-Frequency Mobile Sensor Data

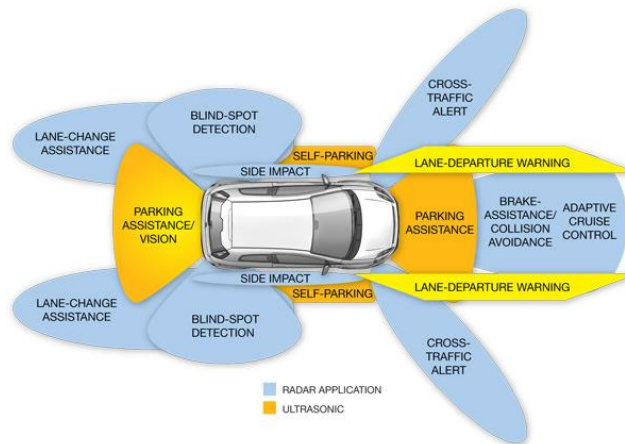
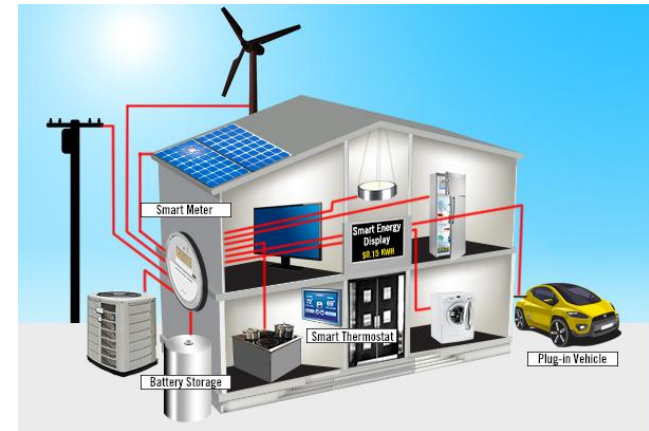


Figure 2 Several driver-assistance systems are currently using radar technology to provide blind-spot detection, parking assistance, collision avoidance, and other driver aids (courtesy Analog Devices).



Health applications are a natural focal point for research using sensor data



Advancing biomedical discovery and improving health through mobile sensor big data

Cornell Tech ♦ Georgia Tech ♦ U. Memphis ♦ Northwestern ♦ Ohio State ♦ Open mHealth
Rice ♦ UCLA ♦ UC San Diego ♦ UC San Francisco ♦ UMass Amherst ♦ U. Michigan ♦ WVU

Enabling mHealth Data Science Research

Traditional Approach

Biomarker data

Hundreds of samples/day

Biomedical Research

Our Proposal

Raw sensor data

Millions of samples/day

Biomarker Identification

Biomarker Development

Biomarker Validation

Biomarker Improvement

Biomedical Research

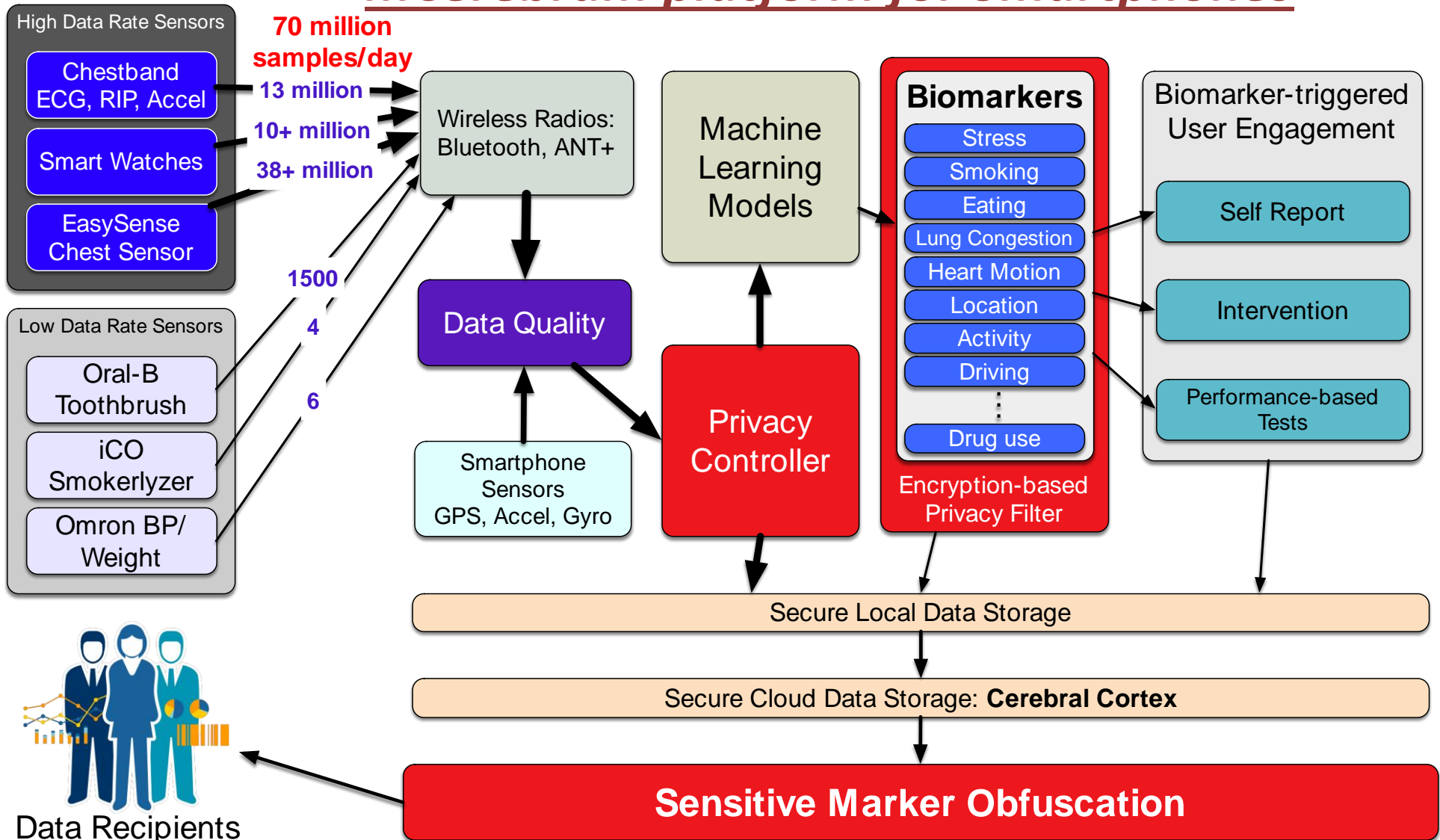


Advancing biomedical discovery and improving health through mobile sensor big data

*Cornell Tech ♦ Georgia Tech ♦ U. Memphis ♦ Northwestern ♦ Ohio State ♦ Open mHealth
Rice ♦ UCLA ♦ UC San Diego ♦ UC San Francisco ♦ UMass Amherst ♦ U. Michigan ♦ WVU*

MD2K Mobile Software Platform (open-source)

mCerebrum platform for Smartphones

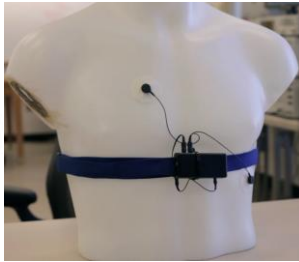


User Studies Using MD2K Software

Study	Target	Participants	Length (days)	Participant-days	Expected samples
Northwestern	Smoking, Stress	210	14	2,940	203 billion
Ohio State	CHF	225	30	6,750	437 billion
UCLA	Oral health	157	180	22,595	878 billion
Vermont	Smoking	90	14	1,260	87 billion
Rice	Smoking	300	30	9,000	622 billion
Utah	Smoking	300	30	9,000	622 billion
Johns Hopkins	Cocaine	25	7	175	11 billion
Totals		1,307	-	51,720	2.86 trillion (~150TB)

Barriers in Conducting Research with High-frequency Mobile Sensor Data

Due to lack of data sharing, everyone needs to collect their own data



Sharing of raw mobile sensor data can accelerate research, but *provenance infrastructure* is needed to enable reproducibility and comparative analysis

Velocity

Hundreds of samples/sec per sensor

Variety

Tens of sensors per sensor

Volume

Gigabytes per day per person

Variability

Variations in attachment, placement, signal quality

Veracity

Multiple biomarkers from same sensor

Validation

Sources of validation for specific biomarkers

mProv: Provenance CI for High-frequency Mobile Sensor Data

- mProv is developing *data models, metadata standards, API's, and runtime support* for annotating sensor data streams with
 - **Source** – sensor type, placement, sampling rate, continuous/episodic
 - **Semantics** – number, probability, class/category;
 - **Provenance** – features and rules applied to obtain a biomarker;
 - **Validation** – specificity, sensitivity, benchmark, gold standard;
 - **Privacy** – user controls exercised and applicable privacy policies
- mProv will enable **replay, interpretability, comparative analysis, and reproducibility**

The mProv Team

Scientific Leadership	Santosh Kumar (PI, U Memphis, analytics); Zachary Ives (Co-PI, U Penn, provenance); Ida Sim (Co-PI, UCSF, metadata); Mani Srivastava (Co-PI, UCLA, privacy)
Scientific Consultants	Emre Ertin (<i>Ohio State</i> , sensor quality); <i>Open mHealth</i> (API integration); James M. Rehg (<i>GA Tech</i> , sensor data analytics)
Advisory Panelists	<p> CISE Advisory Panel: Tanzeem Choudhury (<i>Cornell</i>); Vasant Honavar (<i>Penn State</i>); David Kotz (<i>Dartmouth</i>); Health Advisory Panel: Brian Bot (<i>Sage Bionetworks</i>); Nick Anderson (<i>UC Davis</i>); Industry Advisory Panel: Joe Corkery (<i>Google</i>); Mike O'Reilly (<i>Apple</i>) </p>
Collaborators	Madeleine Ball (Open Humans Project, Participant Recruitment); Gary Wolf (Quantified Self, Participant Recruitment)