

# Collecting **High-frequency** Mobile Sensor Data for **Long-lasting Research Utility**

**Santosh Kumar**

*Director, MD2K Center of Excellence*

*Professor & Lillian and Morrie Moss Chair of Excellence*

Department of Computer Science, University of Memphis



*NIH Big Data to  
Knowledge (BD2K)*



# MD2K Multidisciplinary Team – 20 investigators

## Data Science Research

- Santosh Kumar, *Memphis* (PI)
- Gregory Abowd, Polo Chau, and Jim Rehg, *Georgia Tech*
- Emre Ertin, *Ohio State*
- Deborah Estrin, *Cornell Tech*
- Tyson Condie, Mani Srivastava, *UCLA*
- Deepak Ganesan, Ben Marlin, *UMass*
- Susan Murphy, *Harvard*

## Health Research

- William Abraham, *Ohio State*
- Inbal Nahum-Shani, *Michigan*
- Bonnie Spring, *Northwestern*
- Cho Lam, Dave Wetter, *Utah*
- Vivek Shetty, *UCLA*
- Ida Sim, *UC San Francisco*
- Jaqueline Kerr, *UC San Diego*
- Clay Marsh, *West Virginia*

**Memphis-based headquarter hosts a team of 10 grad students, a postdoc, 3 software engineers, and 6 staff members**



*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*



# Measuring Exposures, Behaviors, and Outcomes

# Mobile Sensors



# Smartwatch



# Chestbands



# Smart Eyeglasses

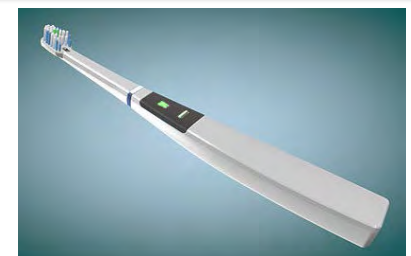
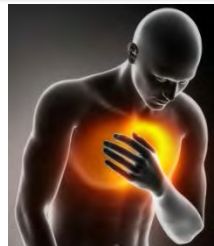
# Exposures



## Behaviors



## Outcomes



*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 Applications – Smoking Cessation & CHF

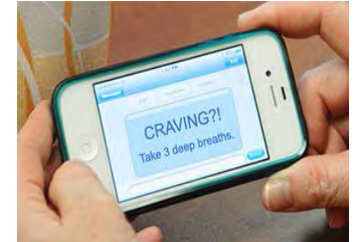
Detect



Predict

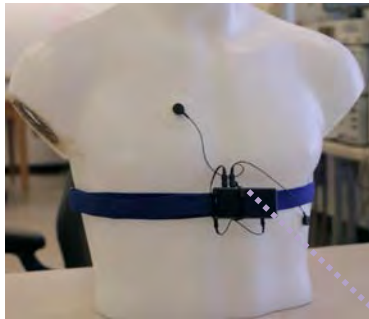


Adapt

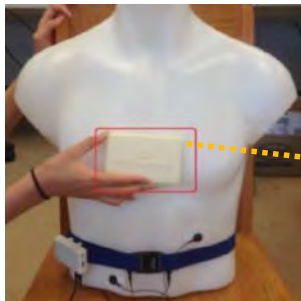




# Mobile Sensor Data Sources in MD2K



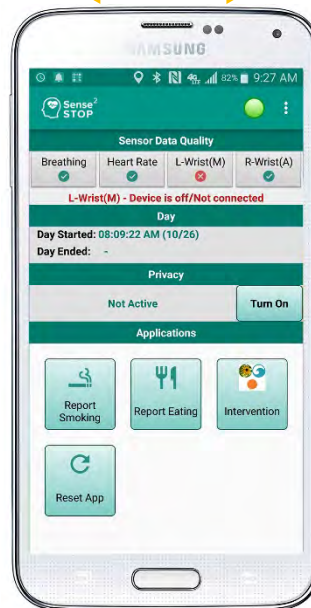
**AutoSense sensors:** ECG, respiration, accelerometers



**EasySense (contactless) sensors:** heart motion, lung motion, lung fluid level



**Smartphone sensors:** GPS, accelerometers, self-report



**Microsoft Band:** accelerometers, gyroscopes, HR

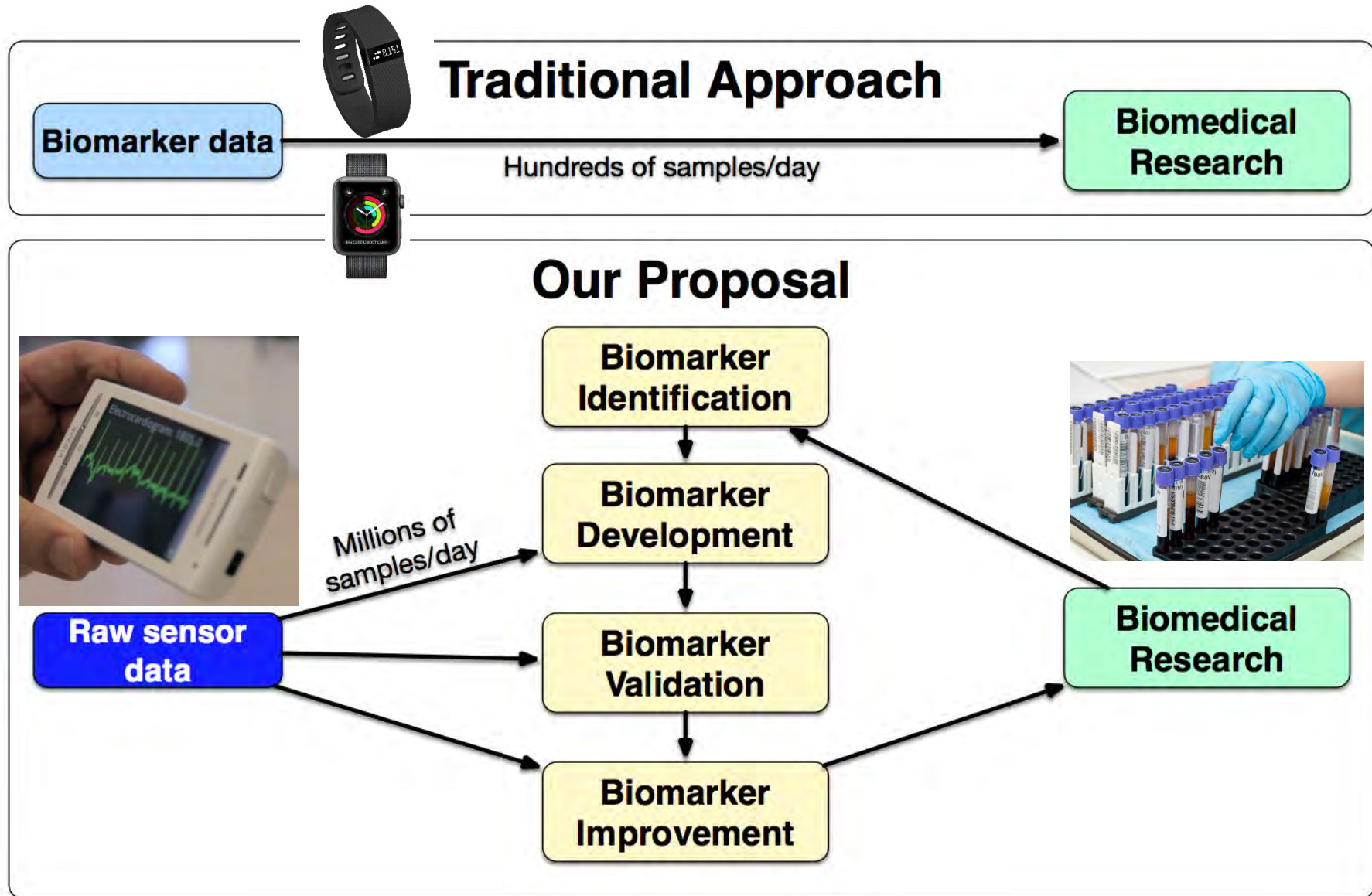


**MotionSense HRV:** accelerometers, gyroscopes, PPG



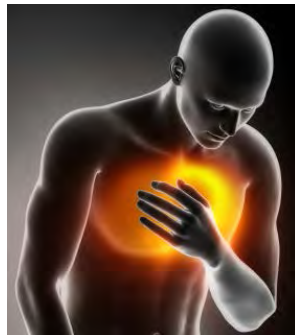
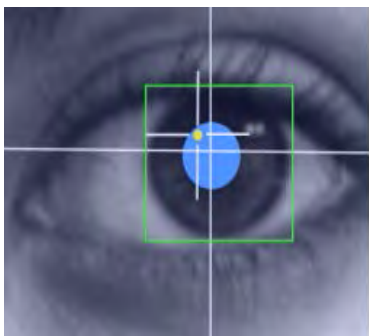
**Smart toothbrush:** brushing, Pressure

# Utility of Collecting High-frequency Sensor Data





# mHealth Biomarkers Developed in MD2K





# Detecting First Lapses in Smoking Cessation

Saleheen, et. al., ACM UbiComp 2015

## Modeling Challenges

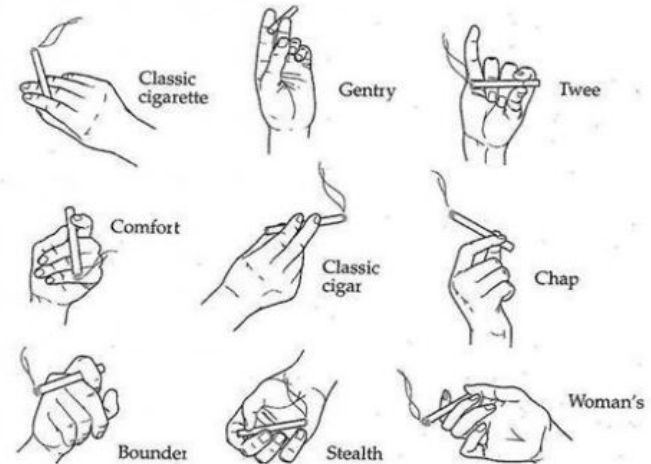
### 1. Ephemeral (very short duration)

- 3~4 sec for each puff
- 10,000 breaths in 10 hours
- 2,000 hand to mouth gestures
- But, only 6~7 positive instances
- Need high recall & low false alarm

### 2. Numerous confounders

- Eating, drinking, yawning

## Wide person & situation variability



<https://www.pinterest.com/pin/566710118890712075/>

## Main Results

- Applied on **smoking cessation data from 61 smokers**
- **Detected 28 (out of 33) first lapses**
- **False alarm rate of 1/6 per day**

## Limitations

- Can't detect if sensor not worn
- Can't detect if data quality is poor
- Needs adaptation for e-cigarettes
- Difficult to validate temporal accuracy of smoking detection

# Sensors-to-Markers-to-Interventions:

## The Case of Sensor-Triggered Stress Intervention

### SENSE

[Ertin, et. al., ACM SenSys'11]



ECG,  
Respiration

Acceleromet  
er

4M samples/day

1M samples/day

### ANALYZE

[Hovsepian, et. al., ACM UbiComp'15]

cStress

Activity

9K samples/day

4K samples/day

### ACT

[Sarker, et. al., ACM CHI'16]

Time series pattern mining

Stress Episode Detection

Intervention Trigger

1-2 Interventions/day

- + High data rate streaming
- + Long battery life
- + High data yield
- + Real-time data quality screening

- + Personalized machine learning models
- + Biomarkers of health, behavior, and environment
- + Validated in lab and field

- + Detect trend in noisy and rapidly varying time series
- + Robust to confounders and data losses
- + Adapt intervention prompts to current context (e.g., driving)



# cStress: Continuous Measure of Stress

Hovsepian, et. al., ACM UbiComp 2015

## 1. A model to convert ECG & respiration sensor data into a continuous measure of stress

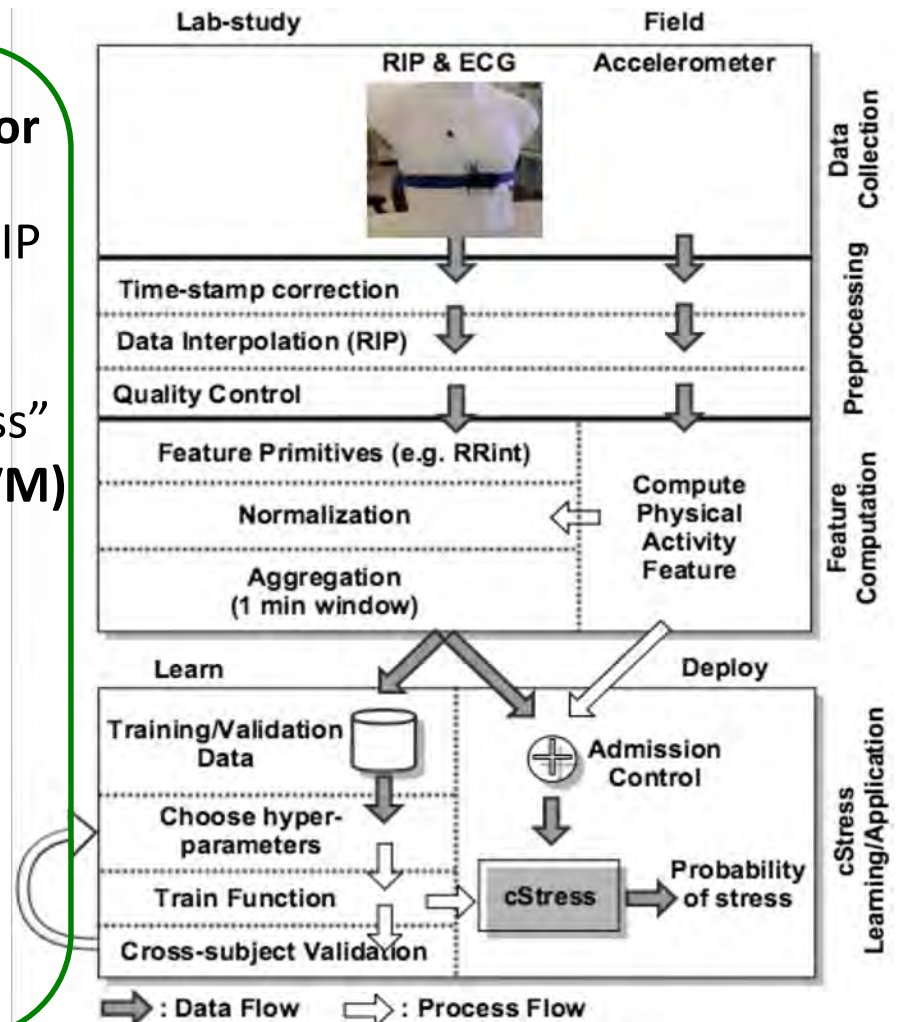
- input: vector of all one-minute ECG + RIP statistical/aggregate features
- output: “stress”/“non-stress” label
- alternative output: probability of “stress”

## 2. Learned with Support Vector Machines (SVM)

- Careful handling of sensor data
- Parameters tuned for optimizing F1 score
- Cross-subject validation for generalizability
- Data loss (0.27 hr/day)

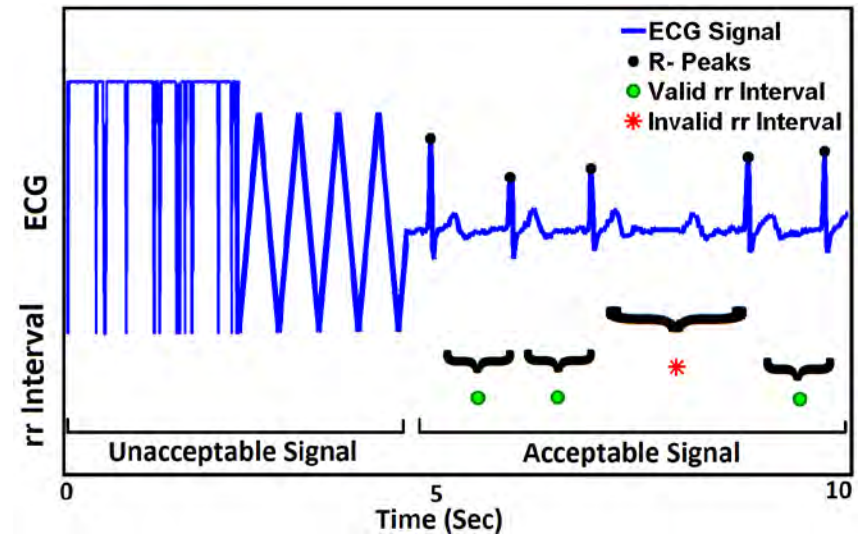
## 3. Validated on independent data sets

- Against lab stress protocol for lab data
- Against self-report for field data



# Feature Computation from ECG

- Screen for data quality
  - Morphology, saturation, etc.
- Find R peaks
  - Pan and Tompkin's algorithm
- Compute R-R interval
- Detect/remove noisy R-R intervals
- Normalize R-R intervals
  - Use “winsorized” (capped) mean/stdev estimates
  - Filter out activity affected data
- Compute one-minute statistics of R-R intervals

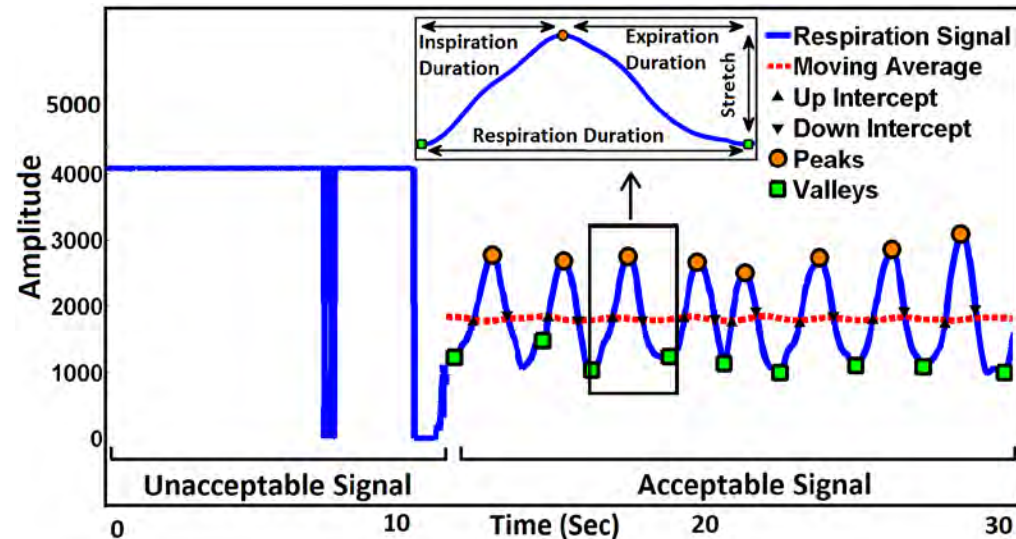


HRV	variance, quartile deviation, low frequency energy (0.1–0.2Hz), medium frequency energy (0.2–0.3Hz), high frequency energy (0.3–0.4Hz), low:high frequency energy ratio
non-HRV	mean, median, 80th percentile, 20th percentile, heart-rate



# Feature Computation from Respiration

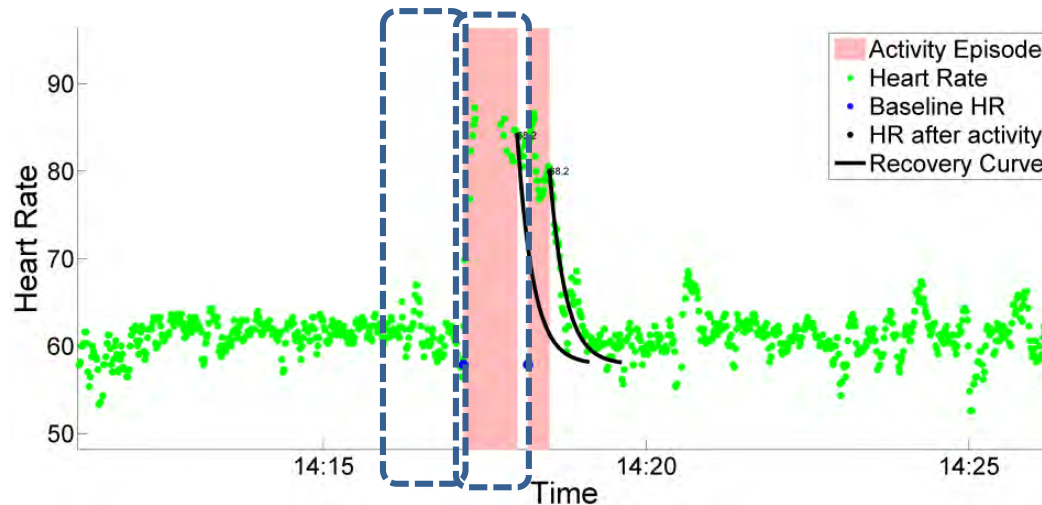
- Screen for data quality
  - Morphology, loosening
- Locate respiration cycles
- Detect/remove invalid cycles
  - Amplitude: > 20% of mean
  - Duration: 0.9 - 12.5 sec
- Compute base features
- Normalize base features
  - As in ECG
- Compute one-minute statistics of base-features



Base Features	Aggregations
inspiration duration, expiration duration, respiration duration, I:E duration ratio, stretch, respiratory sinus arrhythmia (RSA) <sup>1</sup>	mean, median, 80th percentile, quartile deviation
breath-rate <sup>2</sup> , inspiration minute volume <sup>2</sup>	

# Minimizing the Data Loss due to Physical Activity

Sarker, et. al., ACM CHI 2016



Heart-rate  
exponential recovery

$$HRR = HR_{Rest} + (HR_{Peak} - HR_{Rest})e^{-\frac{t-t_0}{\tau}} \quad [\text{Freeman, PCD 2006}]$$

Data Loss due to  
Physical Activity

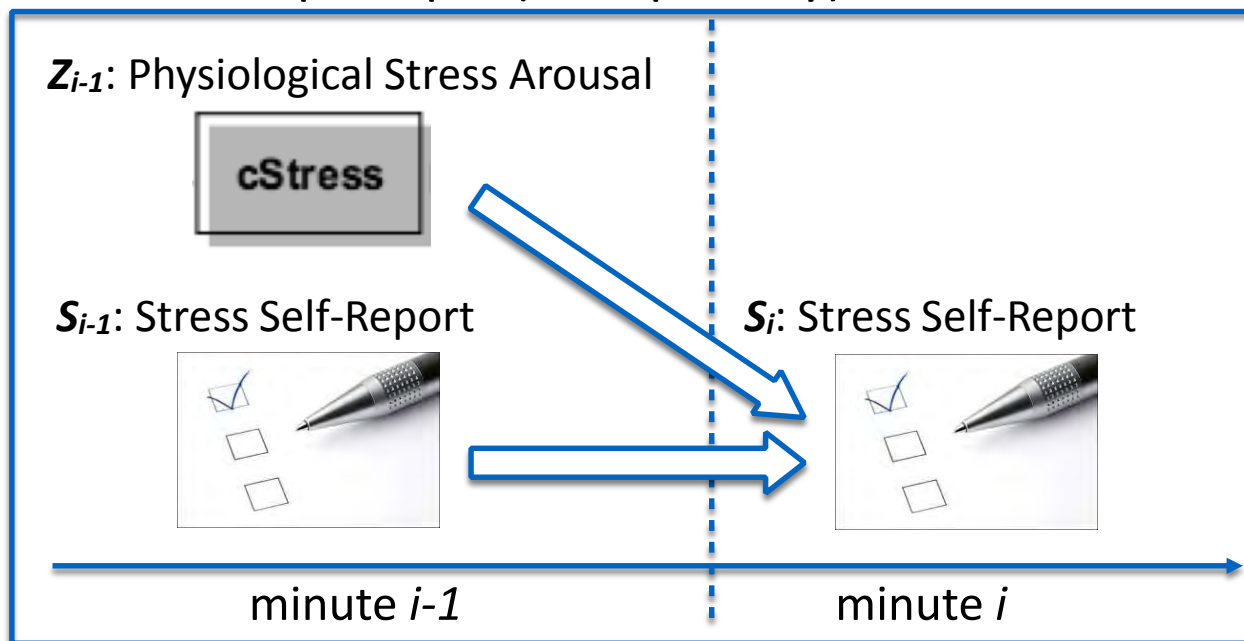
**Recovered: 27.6%**

	Discard 2 minutes	Proposed Method
Activity	22.7%	22.7%
Residue	35.0%	7.4%
Total	57.7%	30.1%



# Training & Validation Methods

- **Lab:** Model trained using lap protocol
  - Public speaking, mental arithmetic, and cold pressor sessions
- **Field:** A Bayesian Network model to map minute-level outputs from cStress to self-reports
  - Random prompts (5-15 per day)



		$S_i$	
$S_{i-1}$	$Z_{i-1}$	1	0
1	1	1	0
1	0	$\alpha$	$1-\alpha$
0	1	$\beta$	$1-\beta$
0	0	0	1

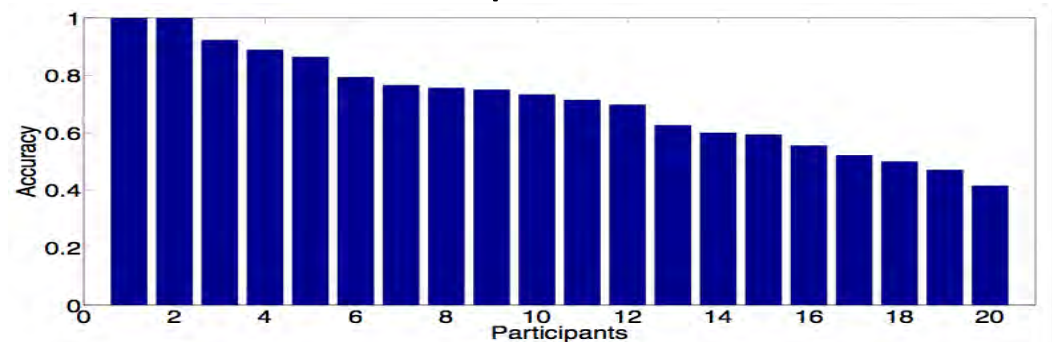
# Validation in Independent Lab and Field Data

- **Lab Validation:** Cross-subject validation with  $n=21$ , 1600 minutes of lab data
- Stress sessions consist of public speaking, mental arithmetic, cold pressor

Feature Set	F1	AUC	Accuracy			C. Kappa	Optimal hyper-parameters		
			Hit-rate	TPR	FPR		$C$	$\gamma$	$bias$
All	0.81	0.96	0.93	0.84	0.05	0.77	90.5097	0.000345267	0.339329
ECG	0.78	0.95	0.92	0.72	0.05	0.73	2	0.00552427	0.340407
HRV	0.56	0.78	0.84	0.55	0.1	0.46	724.077	0.0220971	0.250926
RIP	0.75	0.93	0.90	0.83	0.09	0.69	1448.15	0.000488281	0.308312

- **Field Validation:** 1601 self-report EMA from  $n=23$  over 7 days in the field
- Bayesian network model to map cStress onto self-reported stress data

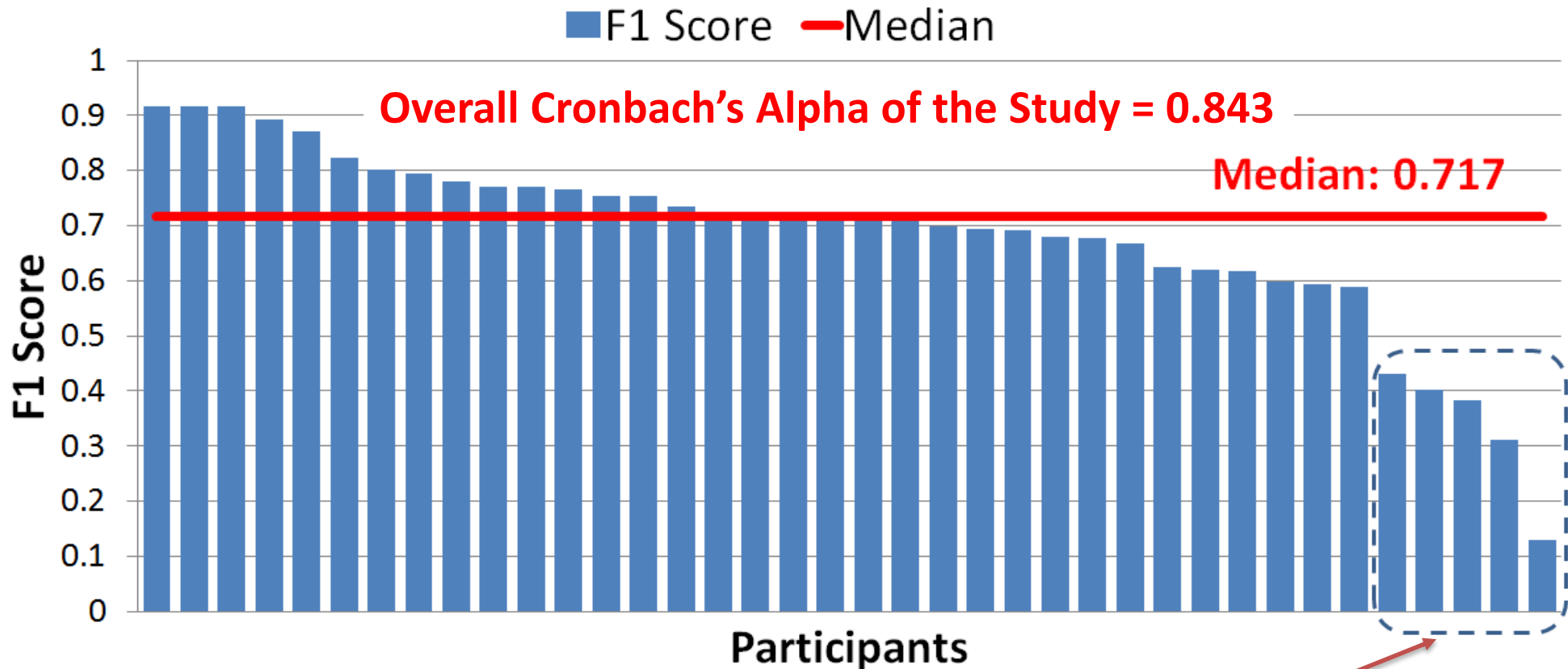
	<i>train</i>	<i>field</i>
Median F1	0.75	0.71
Median AUC	0.85	0.60
Median Accuracy	0.9	0.72





# Validation of cStress on 38 Drug Users Dataset

4 weeks of sensor wearing by polydrug users at NIDA IRP (PI: Dr. Kenzie Preston)

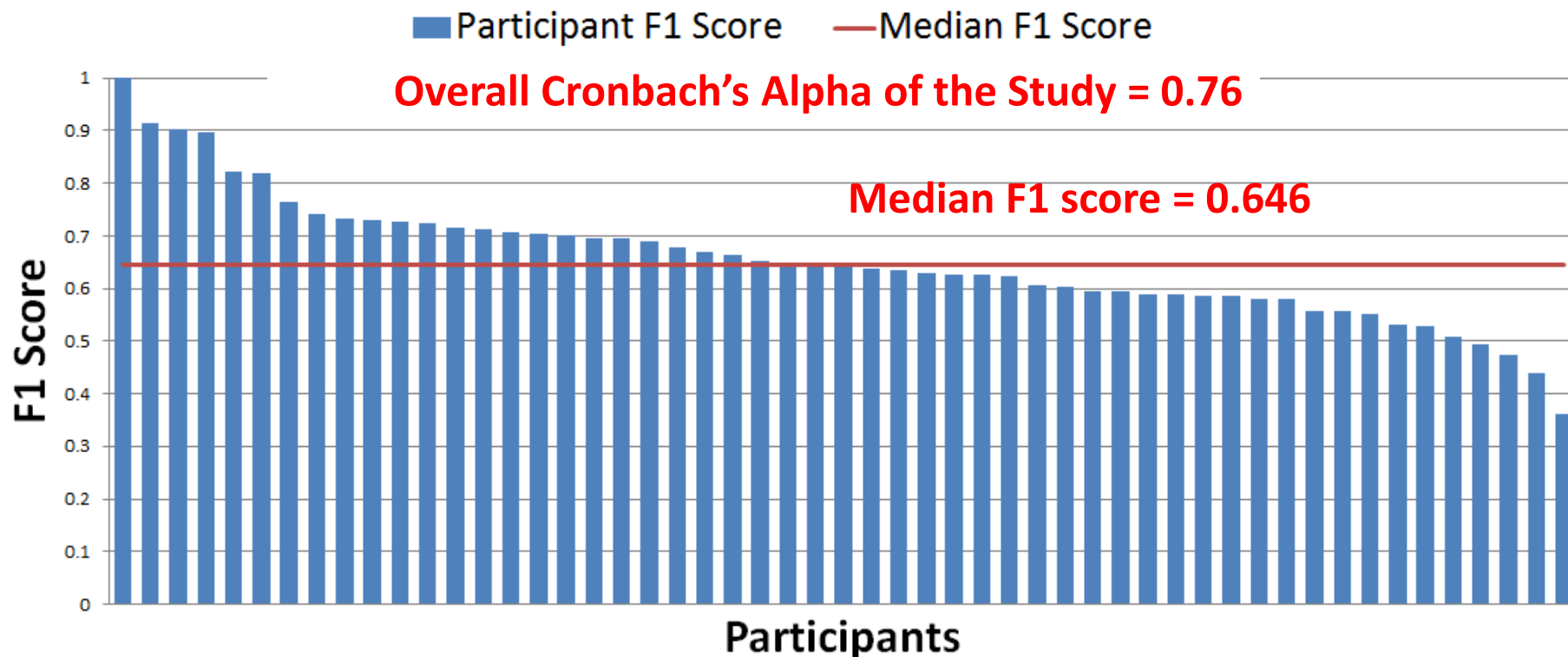


Cronbach's Alpha of last five participants = 0.335

[Sarker, et. al., ACH CHI'16]

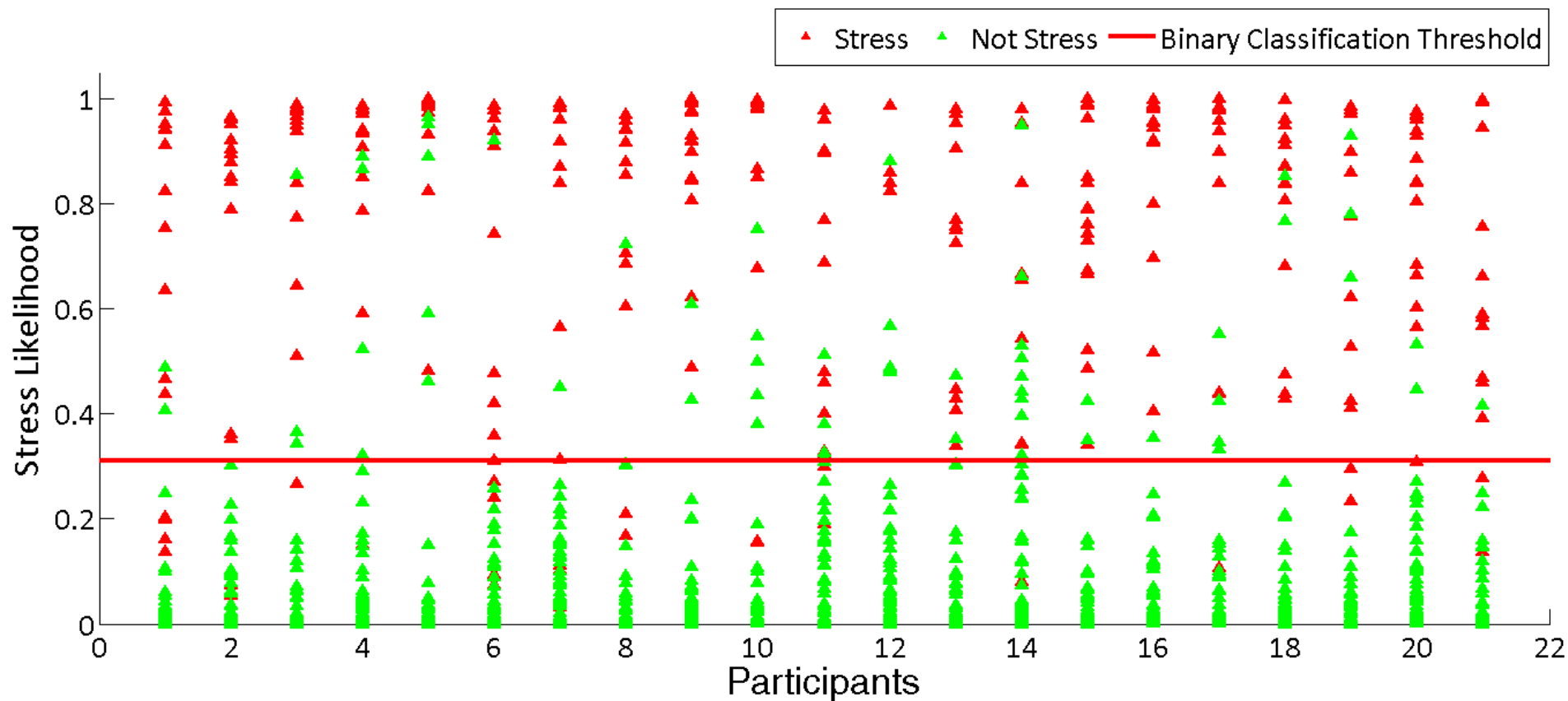
# Validation of cStress on Smoking Data

1 day pre-quit and 3 days post-quit sensor wearing by 61 newly abstinent smokers at UMN (PI: Dr. Mustafa al'Absi)



- **Lower F1 score than other datasets**
  - Imputed missing data, but using simple carry-forward
  - Lower consistency of self-reports (0.76 vs. 0.843)

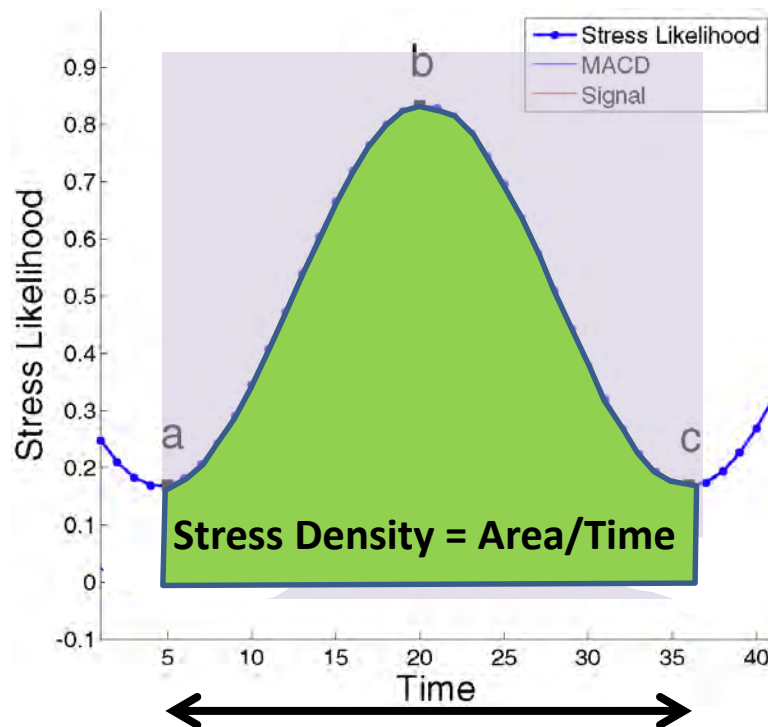
# Stress Likelihood Timeseries



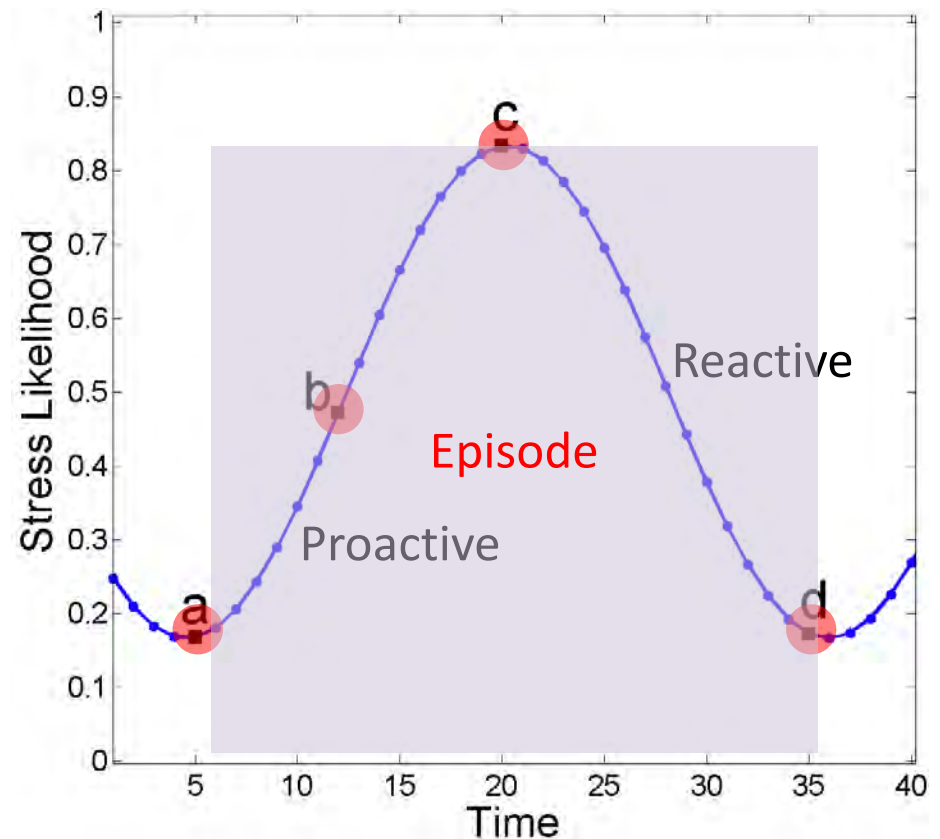


# Mining Stress Episodes in cStress Time Series

Stress Likelihood  $\rightarrow$  Stress Density

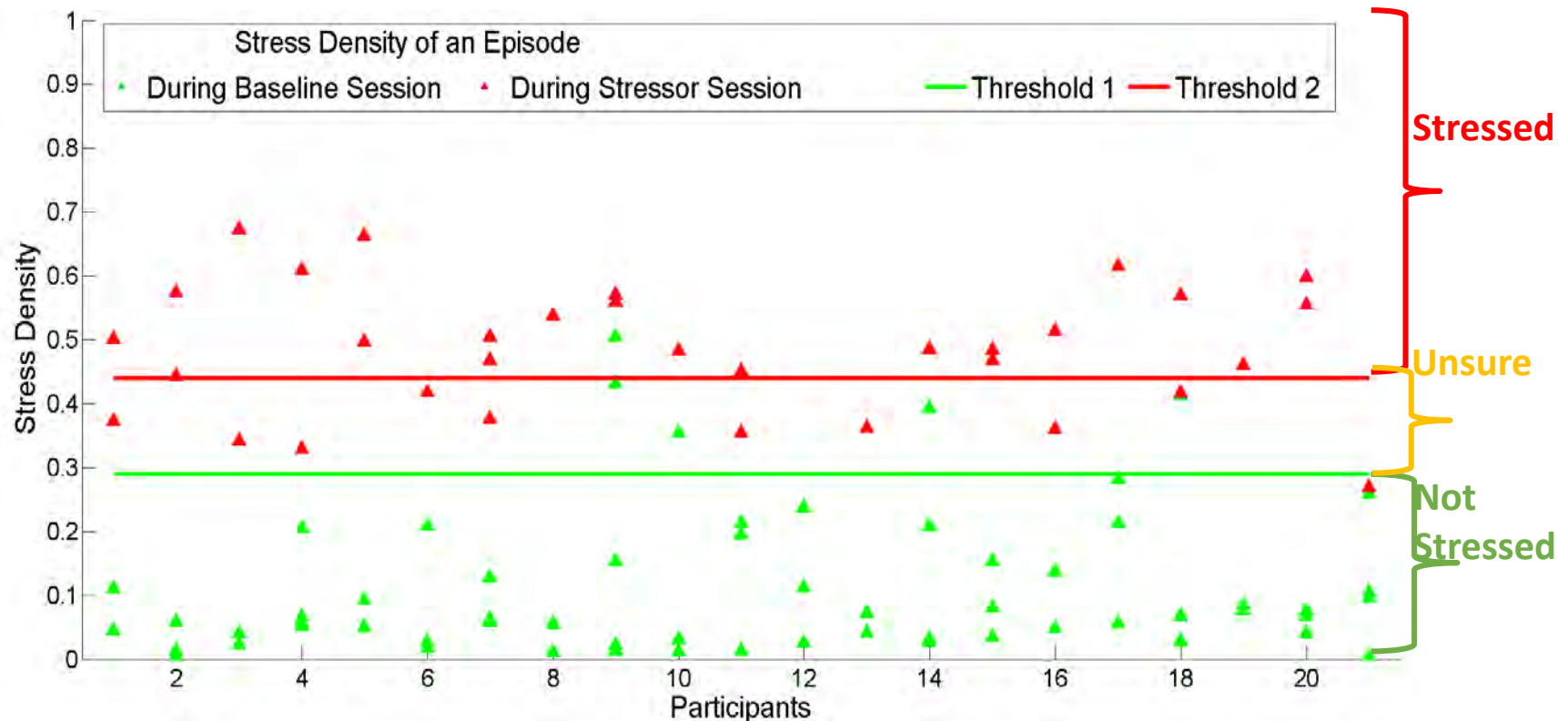


Options for Intervention Timing

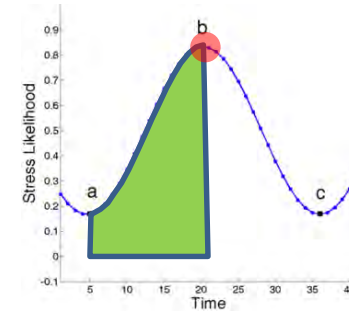
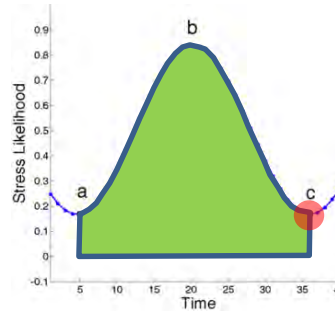


# Generating Intervention Triggers - Model Training

## *Moving from a Single Threshold to Dual Thresholds To Optimize Confidence*



# Thresholds for Reactive Stress Intervention



		Precision and Recall			Precision and Recall	
		95%	90%	85%	90%	85%
Lab Study (Stress Density)	Threshold 1	0.29	0.29	0.29	0.36	0.36
	Threshold 2	0.44	0.42	0.29	0.33	0.36
Field Study (per day)	Not-stress	28.3	28.3	28.3	28.9	29.8
	Unsure	2.7	2.5	0	0.9	0
	Stress	1.5	1.7	4.2	5.1	5.1



# Utility of Collecting High-frequency Sensor Data

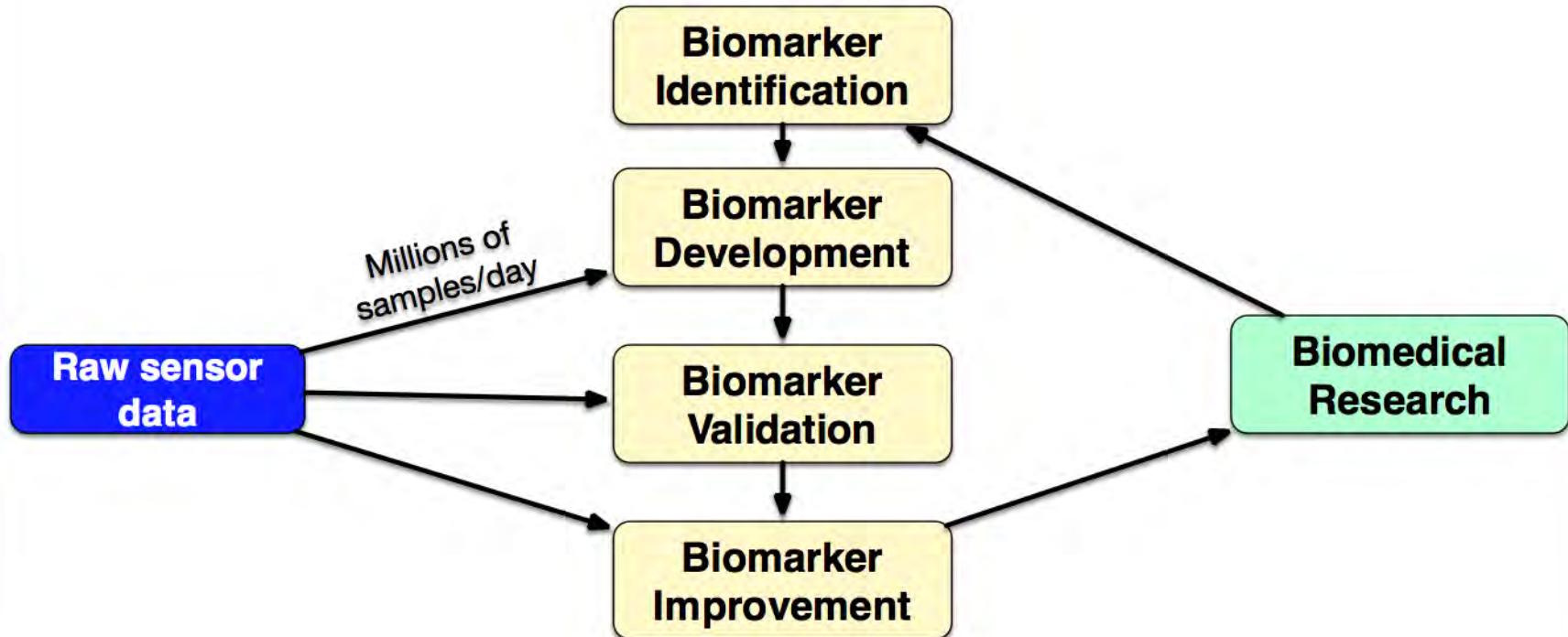
## Traditional Approach

Biomarker data

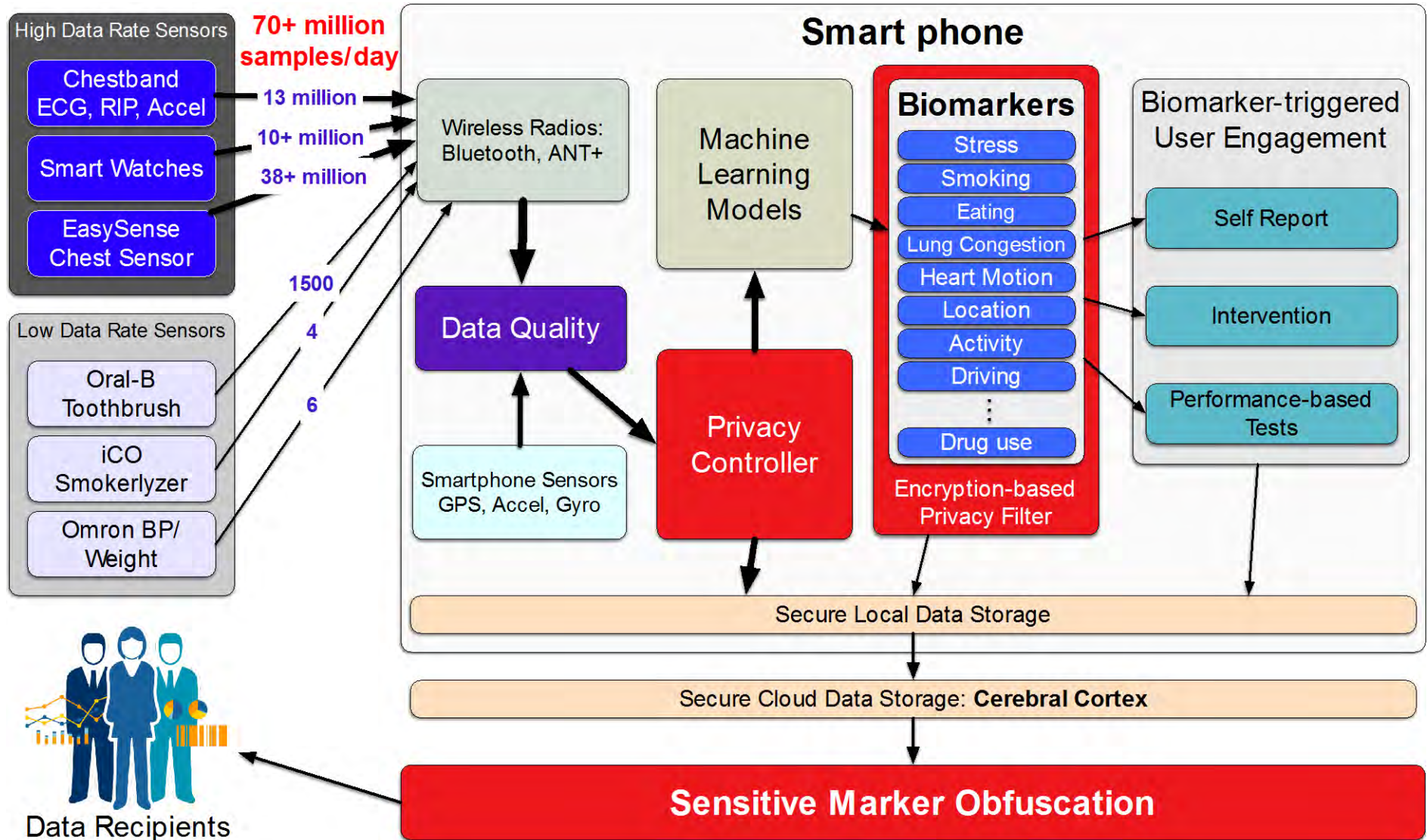
Hundreds of samples/day

Biomedical Research

## Our Proposal



# MD2K Mobile Software Platform (open-source)



# Key Capabilities of mCerebrum

1. Support for high-frequency streaming data
  - 800+ Hz for 70 million samples per day
2. Connectivity to diverse sensors and radio
  - ANT, Bluetooth, Bluetooth Low Energy (BLE), etc.
3. Continuous data collection and real-time data quality monitoring
4. Real-time computation of biomarkers
  - Stress, smoking, driving, activity, etc.
5. Biomarker-triggered notification/intervention



*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



# Field Studies Using MD2K Software

Study	Users	Person-Days	Samples (Billions)
Northwestern (Smoking and Eating)	225	3,150	136
Rice (Smoking)	300	4,200	182
Utah (Smoking)	300	4,200	182
Vermont (Smoking and fMRI)	90	1,260	55
Moffitt (Smoking and Stress)	24	336	15
Ohio State (Heart Failure)	225	6,750	224
UCLA (Oral Health)	162	29,160	968
Johns Hopkins (Cocaine Use)	25	350	18
Dartmouth (Behavior Change)	100	1400	58
Minnesota (Workplace Performance)	800	56,000	2,891
<b>Total</b>	<b>2,251</b>	<b>106,806</b>	<b>4,729</b>

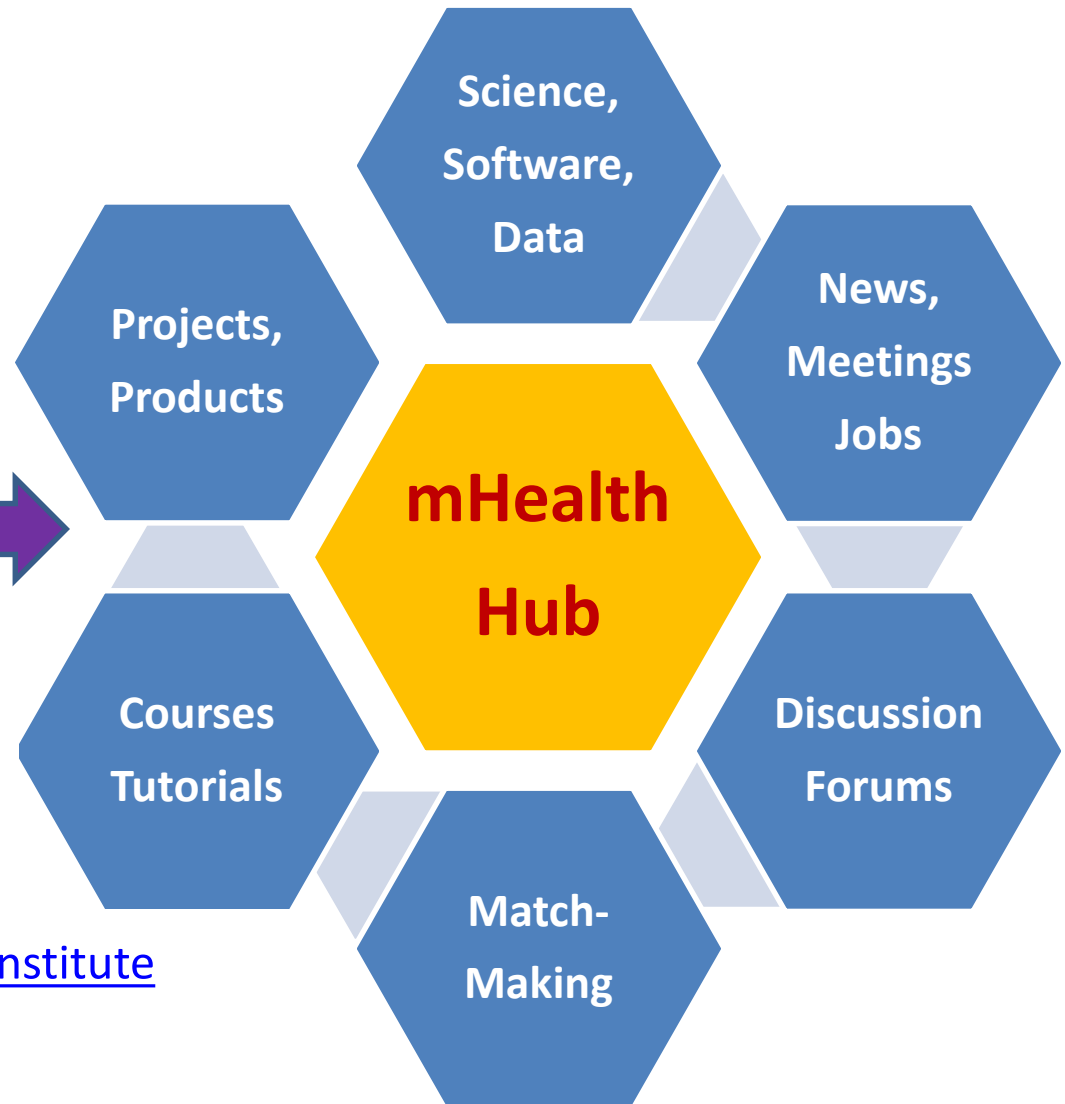
**Entire ecosystem (sensors, software, cloud) to be available end of 2017**

# mHealthHub



**mHealth2017**  
SUMMER TRAINING INSTITUTES

[md2k.org](http://md2k.org)      [info@md2k.org](mailto:info@md2k.org)



27

# mProv: Provenance Cyberinfrastructure for 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



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



# Indicators of Everyday Job Performance



Work hard now. it'll pay off later.



## CORPORATE CITIZENSHIP

### Counterproductive Workplace Behaviour

worktest  
ONLINE ASSESSMENTS



# Healthier, Wealthier, and Happier You

Detect



Predict



Adapt



# For More Information

- MD2K website: [md2k.org](http://md2k.org); Email: [info@md2k.org](mailto:info@md2k.org)
- mProv Website: [mprov.md2k.org/](http://mprov.md2k.org/)
- Software Overview: [md2k.org/software/platform](http://md2k.org/software/platform)
- Software Download: GitHub: [github.com/MD2Korg](https://github.com/MD2Korg)
  - 20+ mCerebrum Android applications
- Software Documentation: [docs.md2k.org](http://docs.md2k.org)
- Questions and Answers: [discuss.md2k.org](http://discuss.md2k.org)
- mHealthHUB: [mhealth.md2k.org/](http://mhealth.md2k.org/)
- mHTI: [mhealth.md2k.org/mhealth-training-institute](http://mhealth.md2k.org/mhealth-training-institute)



*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*