NSF CIF21 DIBBs: mProv: Provenance-Based Data Analytics Cyberinfrastructure for High-Frequency Mobile Sensor Data

Pls: Santosh Kumar, Zachary Ives, Mani Srivastava, Ida Sim Team: John Frommeyer, Timothy Hnat, Nasir Ali, Ananda Tirtha, Sandeep Singh, Nan Zheng, Phillip Hilliard













| torage | Privacy | Query/Replay |
|---------|-----------------------|----------------------|
| and | Provide access | Develop capabilities |
| | control policy and | for querying |
| ly | enforcement | streaming data's |
| | mechanisms for | provenance and for |
| | streaming data, | replaying streams |
| or | using metadata | from archival |
| ource & | annotations | |
| 1 | | |





mProv Use Case: Annotations & Provenance for a Marker Stream Devices run **mCerebrum** and participate in studies Each has sensors with associated metadata Data items stream off each sensor at different rates, and can be *annotated* with provenance and policy information Stream joins integrate the streams from each device, for each owner A stream join integrates the device streams for each owner Based on the provenance annotations, we control access to private data from the home Wrist Sens We accumulate data into **sliding windows** Quality Check and then perform a quality check on recent data

mProv stores, streams, visualizes, & queries *provenance* and annotations for input & marker streams

mProv captures (in a stream and in storage) the output of the pipeline, but also the inputs and steps, and any annotations on them! activity (processing step) *collection ("window")*



For streaming data – we can **query about its provenance and** metadata to, e.g., filter or rank the result

For stored data – we can **merge streams** and **replay precisely**, differentiate among elements, and see if they were produced

