# IEEE Big Data

We had a paper accepted into IEEE Big Data this year. The paper goes over the architecture as well as the various components that make up Brown Dog. If you need to cite Brown Dog this is the paper to use:

S. Padhy, G. Jansen, J. Alameda, E. Black, L. Diesendruck, M. Dietze, P. Kumar, R. Kooper, J. Lee, R. Liu, R. Marciano, L. Marini, D. Mattson, B. Minsker, C. Navarro, M. Slavenas, W. Sullivan, J. Votava, K. McHenry, "Brown Dog: Leveraging Everything Towards Autocuration", *IEEE Big Data*, 2015