# DATA ANALYTICS USING H2O

**F. MOHAMMED**
*Utech LLC location – Madison*

**D. BALTUNOU**
*Utech LLC location – Madison*

**C. DZIK**
*Utech LLC location – Madison*

*Utech LLC, Madison, MS USA*