# APACHE SPARK

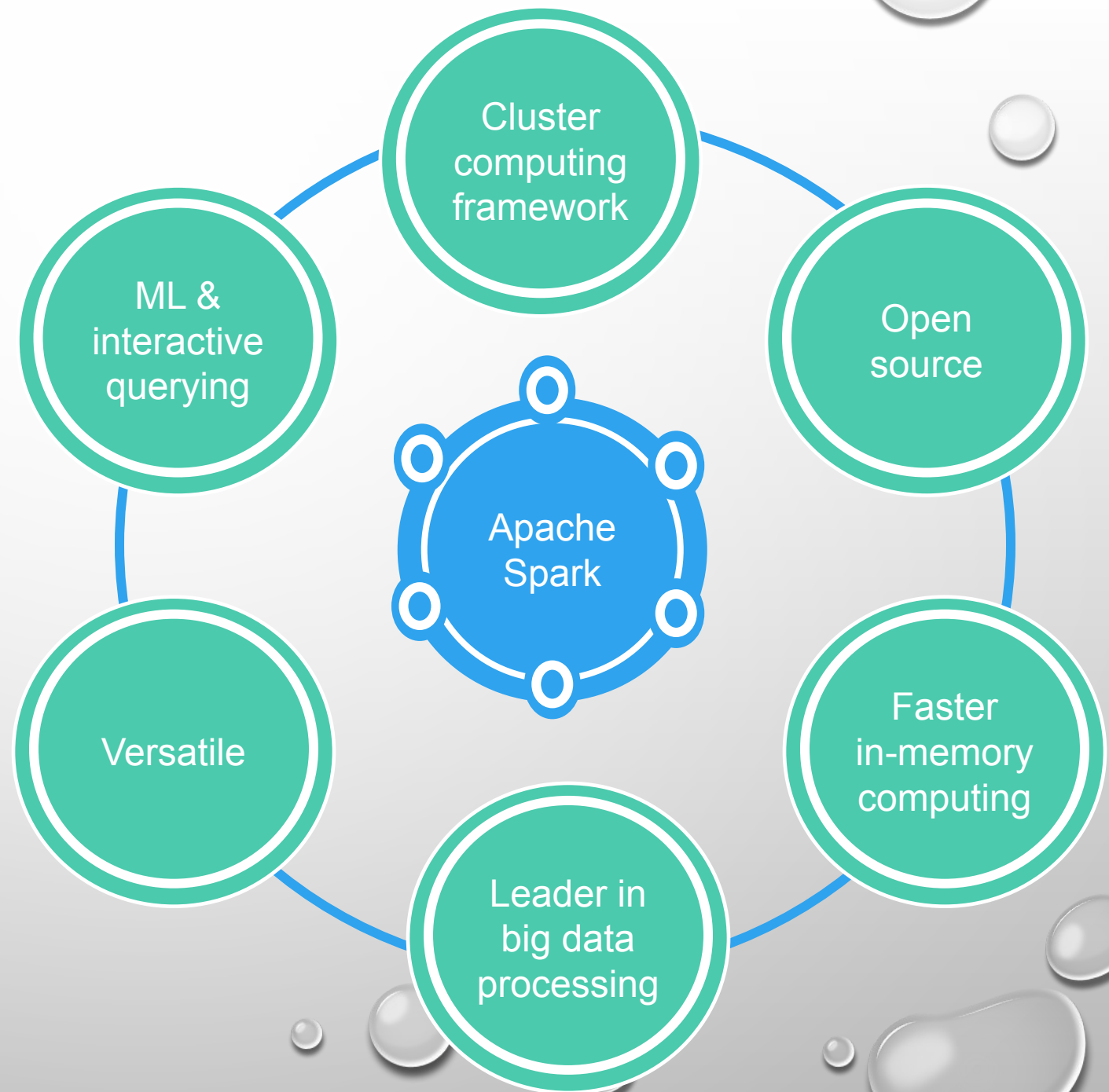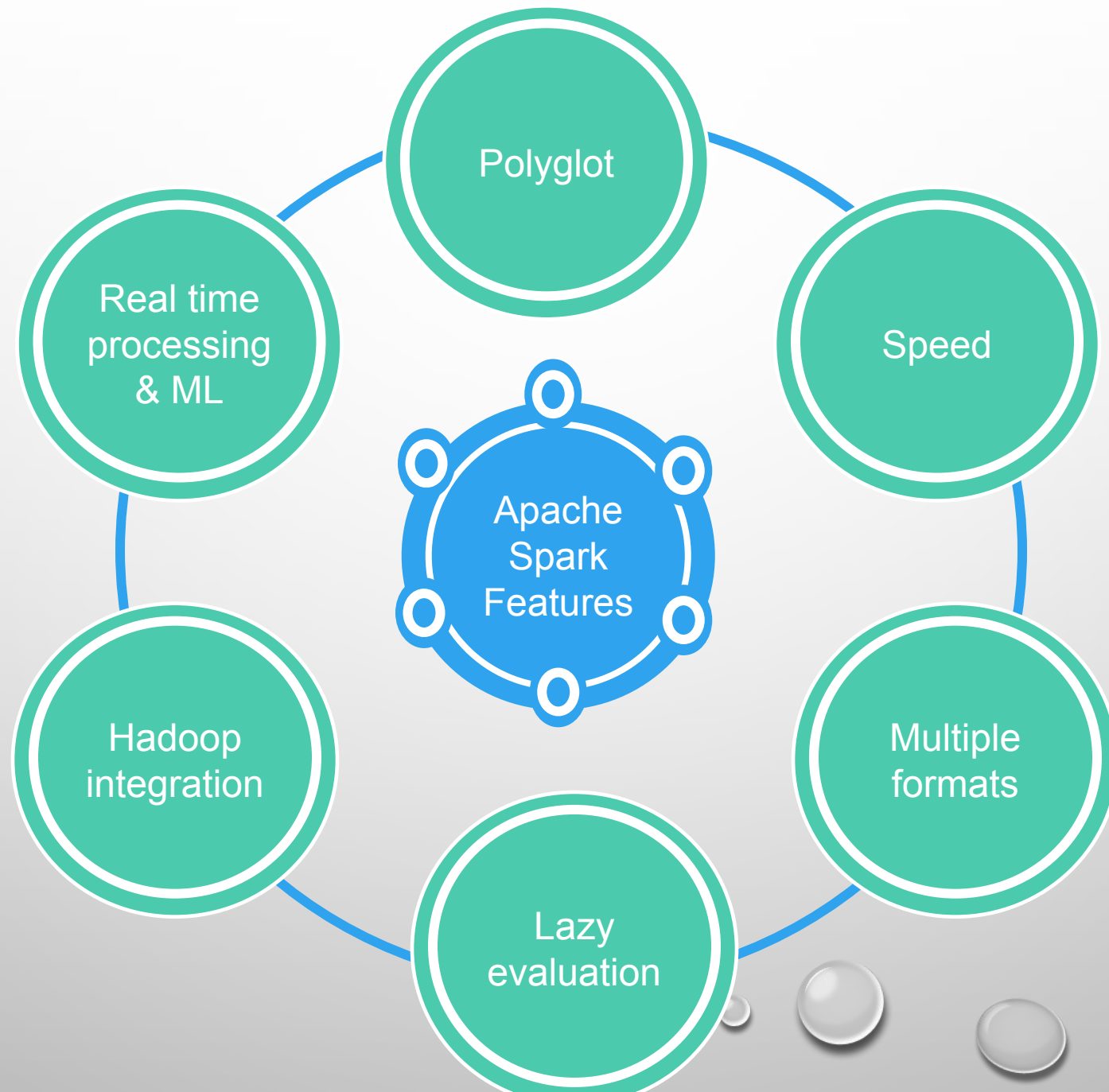# CONTENTS

- APACHE SPARK

- SPARK ARCHITECTURE

- RDD

- DATAFRAME API

- SPARKSQL

- EDA WITH PYSPARK

- PREDICTIVE ANALYSIS WITH SPARK MLIB