(12) **UK Patent Application** (19)**GB** (11)**2610975** (13)**A**

(71) Applicant(s):
    **International Business Machines Corporation
    (Incorporated in USA - New York)
    New Orchard Road, Armonk, New York 10504,
    United States of America**

(72) Inventor(s):
    **Ashish Ranjan
    Arvind Kumar
    Carl Radens**

(74) Agent and/or Address for Service:
    **Elkington and Fife LLP
    Prospect House, 8 Pembroke Road, SEVENOAKS,
    Kent, TN13 1XR, United Kingdom**

(54) Title of the Invention: **Optimal placement of data structures in a hybrid memory based inference computing platform**
    Abstract Title: **Optimal placement of data structures in a hybrid memory based inference computing platform**

(57) Various embodiments are provided for optimized placement of data structures in memory in a computing environment. One or more data structures may be dynamically allocated in a hybrid memory according to weights and activations of the or more data structures in a deep neural network ("DNN"). The hybrid memory includes at least a first memory and a second memory that differ according to write endurance attributes.
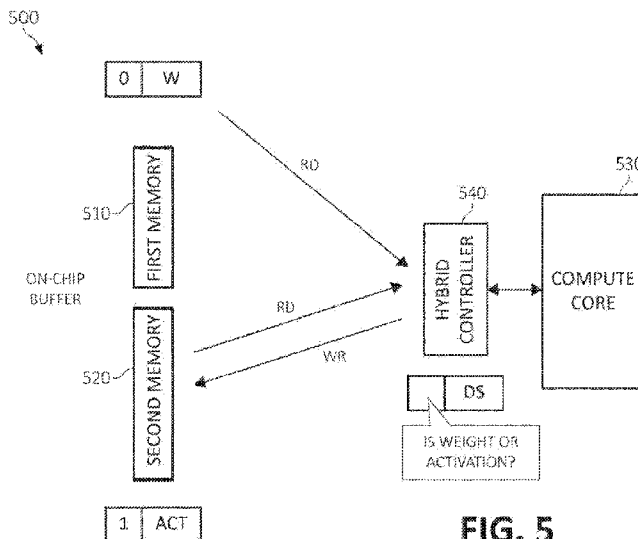


FIG. 5

GB 2610975 A