# iJRASET

## Certificate

It is here by certified that the paper ID : IJRASET63985, entitled

*Open-AI model Efficient Memory Reduce Management for the Large Language Models (LLMs) Serving with Paged Attention of sharing the KV Cashes*

*by*

*Dr. K. Naveen Kumar*

after review is found suitable and has been published in

*Volume 12, Issue VIII, August 2024*

in

*International Journal for Research in Applied Science & Engineering Technology*

*(International Peer Reviewed and Refereed Journal)*

*Good luck for your future endeavors*

**Editor in Chief, iJRASET**

## Certificate

*It is here by certified that the paper ID : IJRASET63985, entitled*

*Open-AI model Efficient Memory Reduce Management for the Large Language Models (LLMs) Serving with Paged Attention of sharing the KV Cashes*

*by*

*Mr. Sreedhar Ambala*

*after review is found suitable and has been published in*

*Volume 12, Issue VIII, August 2024*

*in*

*International Journal for Research in Applied Science & Engineering Technology*

*(International Peer Reviewed and Refereed Journal)*

*Good luck for your future endeavors*

**Editor in Chief, iJRASET**

*Certificate*

_It is here by certified that the paper ID : IJRASET63985, entitled_

_Open-AI model Efficient Memory Reduce Management for the Large Language Models (LLMs) Serving with Paged Attention of sharing the KV Cashes_

_by_

_Dr. M. B Raju_

_after review is found suitable and has been published in_

_Volume 12, Issue VIII, August 2024_

_in_

International Journal for Research in Applied Science & Engineering Technology

(International Peer Reviewed and Refereed Journal)

Good luck for your future endeavors

**Editor in Chief, iJRASET**

# iJRASET

**International Journal for Research in Applied Science & Engineering Technology**

## Certificate

_It is here by certified that the paper ID : IJRASET63985, entitled_

_**Open-AI model Efficient Memory Reduce Management for the Large Language Models (LLMs) Serving with Paged Attention of sharing the KV Cashes**_

_by_

_**Dr. Sai Hareesh Anamandra**_

_after review is found suitable and has been published in_

_Volume 12, Issue VIII, August 2024_

_in_

**International Journal for Research in Applied Science & Engineering Technology**

_(International Peer Reviewed and Refereed Journal)_

_Good luck for your future endeavors_

**Editor in Chief, iJRASET**

# iJRASET







ISRA Journal Impact Factor: **7.429**













*It is here by certified that the paper ID : IJRASET63985, entitled*

*Open-AI model Efficient Memory Reduce Management for the Large Language Models (LLMs) Serving with Paged Attention of sharing the KV Cashes*

*by*

*Mr. Dhanunjaya Rao Kodali*

*after review is found suitable and has been published in*
*Volume 12, Issue VIII, August 2024*
*in*

*International Journal for Research in Applied Science & Engineering Technology*
*(International Peer Reviewed and Refereed Journal)*
*Good luck for your future endeavors*

Editor in Chief, iJRASET