# Comparative Genomics Analyses of Selected Microbes with Open Pan-genomes Highlight Their Evolutionary Dynamics at Functional Level

A Thesis

*to be submitted by*

**Vikas Sharma (Roll No D12090)**

*for the award of the degree*

*of*

**Doctor of Philosophy**



*School of Basic Sciences*

**Indian Institute of Technology Mandi**

**Kamand, Himachal Pradesh-175005, India**

February, 2019

# Declaration by the Scholar

I hereby declare that the entire work embodied in this Thesis is the result of investigations carried out by me in the **School of Basic Sciences**, Indian Institute of Technology Mandi, India, under the supervision of **Dr. Tulika Prakash Srivastava**, and that it has not been submitted elsewhere for any degree or diploma. In keeping with the general practice, due acknowledgements have been made wherever the work described is based on finding of other investigators.

Place:                                    Signature

Date:                                     Name: **Vikas Sharma**

# <u>Declaration by the Research Advisor</u>

I hereby certify that the entire work in this Thesis has been carried out by **Vikas Sharma** under my supervision in **School of Basic Sciences**, Indian Institute of Technology Mandi, and that no part of it has been submitted elsewhere for any Degree or Diploma.

Signature:

Name of the Guide: Dr. Tulika Prakash Srivastava

Date:

*Dedicated to My Mother*

# Acknowledgement

# Table of Contents