# World-Wide Virtual Machine:

## A Metacomputing Environment

## Integrating World-Wide Web and High Performance

## Computing and Communications Technologies

by

KIVANC DINCER

B.S., Middle East Technical University, Turkey, 1989
M.S., Iowa State University, 1991

DISSERTATION

Submitted in partial fulfillment of the requirements for the
degree of Doctor of Philosophy in Computer Science
in the Graduate School of Syracuse University

August 1997

Approved ————————————————
Professor Geoffrey C. Fox

Date ————————————————

# Abstract

This thesis discusses the major issues in building a metacomputing environment based on World-Wide Web (WWW) and High Performance Computing and Communication (HPCC) technologies and describes the design and implementation of such an environment called World Wide Virtual Machine (WWVM). The presented work helps to carry much of the past decade's work in HPCC technologies to the larger WWW domain.

The WWVM exploits the open interfaces brought by Web servers. It extends the servers via Common Gateway Interface (CGI) extensions and uses PVM daemons and low-level protocols such as HTTP, TCP/IP, and UDP/IP in order to combine remote networked computers as a single machine. The WWVM can work in stand-alone, message-passing, and dataflow modes and provides the interoperability of many diverse software and hardware components. The stand-alone mode allows computations to be performed on a remote WWVM server or any other machine coordinated by a server. In the message-passing mode, the WWVM is capable of executing message-passing PVM and MPI programs, and High Performance Fortran (HPF) programs compiled by the Syracuse Fortran 90D/HPF compiler, as well as parallel programs using PCRC or Global Arrays runtime support libraries. In the dataflow mode, WWVM's coordination layer interprets a simple task flow (data-dependency) description language to deduct the dataflow patterns between different WWVM nodes.

The WWVM supplies an integrated, Web-based programming environment and gives pervasive access to remote WWVM facilities from any platform (Unix, PC, or Mac) using a standard Web browser. Client-side Web technologies such as HTML, JavaScript, plug-ins, and Java supply a platform-independent graphical user interface and visualization capabilities that

include analyzing data output from programs and performance information recorded in Pablo's SDDF format. The associated data wrapper libraries provide real-time, application-specific data

programs and client-side Java applets.

# Table of Contents

# List of Tables

# List of Figures

# Acknowledgements

I appreciate my advisor, Professor Geoffrey C. Fox, for his guidance throughout my research work and for his wise and acute observations on how to improve my work. I gratefully acknowledge Dr. Anne Trefethen from Cornell Theory Center, Dr. Simon Catterall, Dr. Wojtek Furmanski, Dr. Xiaoming Li, and Dr. Stephen Taylor for serving on my defense committee.

I would like to extend my deepest thanks to Dr. Alok Choudhary, Dr. Tomasz Haupt, Dr. Sanjay Ranka, Dr. David Bernholdt, and Mr. Donald Leskiw with whom I worked and learned from during the last four years.

Many thanks are also due Dr. Salim Hariri, Dr. Xiaoming Li's students from China, and a number of anonymous reviewers who read earlier technical paper versions of this manuscript and provided useful comments. Their excellent suggestions greatly helped me in preparing the final revision.

My sincere thanks go to Ms. Kathy Barbieri and Ms. Caroline Hecht from Cornell Theory Center for providing valuable suggestions for my work.

I am especially grateful to Mrs. Elaine Weinman for her help in English and for her patience in proofreading endless revisions of this manuscript.

I would like to express my appreciation to many friends who contributed to my work with their encouragement. Among them, Mr. Haluk Topcuoglu deserves special thanks. His moral support and suggestions were of great value.