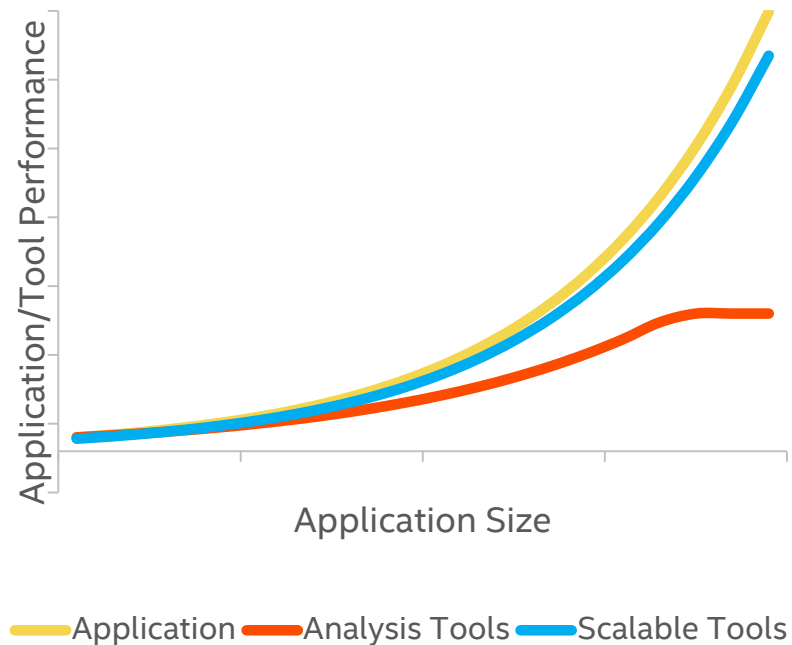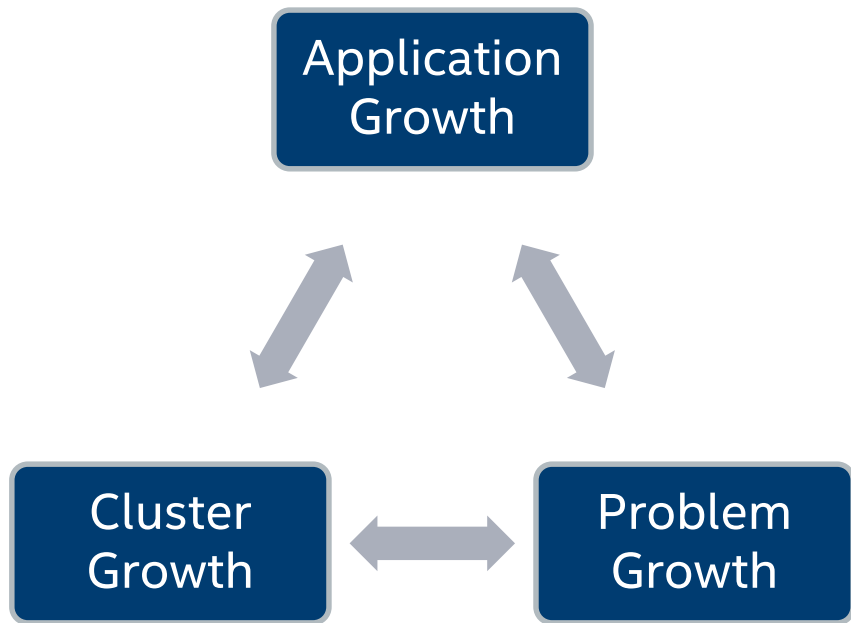# Agenda

Scaling with MPI Performance Snapshot

Tuning MPI Performance with Intel® Trace Analyzer and Collector

MPI SCALING ANALYSIS CHALLENGES

# MPI Scaling Analysis Challenges
# To Exascale… and Beyond



**Application Growth**

**Cluster Growth**

**Problem Growth**

Application/Tool Performance

Application Size

Application — Analysis Tools — Scalable Tools

# MPI PERFORMANCE SNAPSHOT (MPS)

# Why MPI Performance Snapshot (MPS)?

- Advantages
  - Get an initial profile of the application very quickly
  - Performance variation at scale can be detected and triaged quickly
  - Provides development recommendations to developers based on analysis
    - Intel® Trace Analyzer and Collector or Intel® VTune™ Amplifier XE for deeper analysis
  - Easy to use out of the box functionality

- Benefits
  - Difficult performance issues are easier to spot
  - Application performance guidance is obtained easily
  - Experienced & non-experienced developers can adopt quickly

# 2 MPI Performance Snapshots

- Bundled version:
  - is a part Parallel Studio Cluster Edition
  - relies on Intel VTune Amplifier
  - Intel MPI supports MPS (Hydra knows some environemt variables)

- Standalone version:
  - Can be downloaded with no charge.
  - Contains Application Perfomance Snapshot (APS) to collect Hardware counters
  - Has got a launcher script

# What's new

Absolutely new collector which can collect information from all MPI functions.

Collector can work with any MPICH-based MPI implementation (OpenMPI is not supported so far)

Only one library to collect statistics and all other metrics (except Gflops, CPI, Memory bound)

Collector uses only MPI standard calls to support compatibility

New mechanism of MPI imbalance collection based on MPI_T_ mechanism – iMPI only

rdpmc timer for really low intrusion

Binary file format for MPI statistics

Each file writes to its own file

Different levels of statistics (MPS_STAT_LEVEL=1…5) 5 is default now.

New parser for binary statistics. (Still supports native iMPI statistics)

Absolutely new HTML report. All metrics on one page.

HTML report now shows top 5 MPI functions

# What's not supported

MPS doesn't print statistics on finalization

PAPI library – we cannot rely on this library

MPI_Pcontrol() is not supported yet.

No MPI imbalance in other MPI implementations (might be added later)

OpenMP imbalance can be caught from Intel OpenMP library only.

MPI functions should not be called from OpenMP regions.

VTune Amplifier requires RedHat 6 and later and `perf` utility supported by kernel

# MPS USAGE

# How to run

"Bundle"

    $ mpirun –mps –n N app_name

Standalone

    $ source mpsvars.sh --vtune

    $ mpirun –n N mpsrun.sh app_name

`mpsrun.sh` is a script which sets needed environment variables including LD_PRELOAD.

In VTune mode it runs `amplxe-cl` for each process

# MPS Output

## Summary

## Files and folders:

- stats.txt
  - MPI statistics

- app_stat.txt
  - MPS collector statistics

- _mps/results.<node>/
  - VTune results

```
================== GENERAL STATISTICS ==================
Total time:    448.391 sec (All ranks)
        MPI:      40.73%
    NON_MPI:      59.27%

WallClock :
        MIN :            89.594 sec (rank 1)
        MAX :            89.975 sec (rank 4)

================ MEMORY USAGE STATISTICS ================
All ranks:    226.969 MB
        MIN:      24.172 MB (rank 2)
        MAX:      96.465 MB (rank 0)

================ MPI IMBALANCE STATISTICS ================
MPI Imbalance:            31.798 sec          7.092% (All ranks)
            MIN:           2.219 sec          2.467% (rank 4)
            MAX:           9.157 sec         10.219% (rank 0)
```

# HTML Reporting

**MPI Performance Snapshot** Summary



WallClock time: 21.78 sec
Total application lifetime. The time is elapsed time for the slowest process. This metric includes the MPI Time and the Computation time below.

MPI Time: 17.28 sec — 80.10%
Time spent inside the MPI library. High values are usually bad.
This value is HIGH. The application is Communication-bound. More details...

MPI Imbalance: 7.44 sec — 34.47%
Mean unproductive wait time per process spent in the MPI library calls when a process is waiting for data. This time is part of the MPI time above. High values are usually bad.
This value is HIGH. The application workload is NOT well balanced between MPI ranks. More details...

Computation Time: 4.29 sec — 19.90%
Mean time per process spent in the application code. This is the sum of the OpenMP Time and the Serial time. High values are usually good.
This value is LOW.

OpenMP Time: 19.46 sec — 90.19%
Mean time per process spent in the OpenMP parallel regions. High values are usually good and indicate that the application is well-threaded.
This value is HIGH.

OpenMP Imbalance: 17.05 sec — 79.03%
Mean unproductive wait time per process spent in OpenMP parallel regions (normally at synchronization barriers). High values are usually bad.
This value is HIGH. The application's OpenMP work sharing is NOT well load-balanced. More details...

Serial Time: 0.00 sec — 0.00%
Mean application time per process spent outside OpenMP parallel regions. High values may be good or bad depending on the application algorithm.
This value is NEGLIGIBLE. This application is well parallelized via OpenMP directives.

WallClock time: 21.78 sec

MPI Time: 17.28 sec — 80.10%
MPI Imbalance: 7.44 sec — 34.47%
Computation Time: 4.29 sec — 19.90%
OpenMP Time: 19.46 sec — 90.19%
OpenMP Imbalance: 17.05 sec — 79.03%
Serial Time: 0.00 sec — 0.00%

Application: build/heart_demo
Number of ranks: 17
Used statistics: app_stat_20160310-035458.txt, stats_20160310-035458.txt
Creation date: 2016-03-10 03:55:21

# MPS HTML Report Breakdown – MPI Time

MPI Time – Time spent in MPI calls

MPI Imbalance – MPI time spent waiting

Lower is better

If MPI Time or MPI Imbalance are high, use Intel® Trace Analyzer and Collector to investigate and optimize MPI usage



```
■ MPI Time: 17.28 sec                                                        80.10%
Time spent inside the MPI library. High values are usually bad.
This value is HIGH. The application is Communication-bound. More details...

   ▨ MPI Imbalance: 7.44 sec                                                 34.47%
   Mean unproductive wait time per process spent in the MPI library calls when a process is waiting for data. This time is part
   of the MPI time above. High values are usually bad.
   This value is HIGH. The application workload is NOT well balanced between MPI ranks. More details...
```

# MPS HTML Report Breakdown – OpenMP Time

OpenMP Time – Computation time spent in OpenMP parallel regions – higher is better

OpenMP Imbalance – OpenMP Time spent waiting – lower is better

If OpenMP Imbalance is high – recommend using Intel® VTune™ Amplifier XE

If OpenMP Time is low – Intel® Advisor to find opportunities to add more threading



OpenMP Time: 19.46 sec      90.19%
Mean time per process spent in the OpenMP parallel regions. High values are usually good and indicate that the application is well-threaded.
This value is HIGH.

OpenMP Imbalance: 17.05 sec      79.03%
Mean unproductive wait time per process spent in OpenMP parallel regions (normally at synchronization barriers). High values are usually bad.
This value is HIGH. The application's OpenMP work sharing is NOT well load-balanced. More details…

# TUNING MPI APPLICATION PERFORMANCE WITH INTEL® TRACE ANALYZER AND COLLECTOR

# Intel® Trace Analyzer and Collector

**Value Proposition**

| | |
|---|---|
| **What** | • **Intel's High Performance MPI Communications Profiler & Analyzer for Scalable HPC Development** |
| **Why** | • **Scale Performance –** Perform on More Nodes<br>• **Scale Forward –** Multicore and Manycore Ready<br>• **Scale Efficiently –** Tune & Debug on More Nodes |
| **How** | • **Visualize  -** Understand parallel application behavior<br>• **Evaluate -** Profiling statistics and load balancing<br>• **Analyze –** Automated analysis of common MPI issues<br>• **Identify –** Communication hotspots |

# Intel® Trace Analyzer and Collector Overview

Intel® Trace Analyzer and Collector helps the developer:

- Visualize and understand parallel application behavior
- Evaluate profiling statistics and load balancing
- Identify communication hotspots

## Features

- Event-based approach
- Low overhead
- Excellent scalability
- Powerful aggregation and filtering functions
- Performance Assistance and Imbalance Tuning
- NEW in 9.1: MPI Performance Snapshot

**API and *-tcollect***

***-trace***

**Intel® Trace Collector**

**Trace File (.stf)**

**Intel® Trace Analyzer**

**Source Code**

**Compiler**

**Objects**

**Linker**

**Binary**

**Runtime**

**Output**

# Strengths of Event-based Tracing

| | |
|---|---|
| **Predict** | Detailed MPI program behavior |
| **Record** | Exact sequence of program states – keep timing consistent |
| **Collect** | Collect information about exchange of messages: at what times and in which order |

An event-based approach is able to detect temporal dependencies!

# Multiple Methods for Data Collection

| Collection Mechanism | Advantages | Disadvantages |
|---|---|---|
| **Run with –trace or preload trace collector library.** | **Automatically collects all MPI calls, requires no modification to source, compile, or link.** | **No user code collection.** |
| Link with –trace. | Automatically collects all MPI calls. | No user code collection. Must be done at link time. |
| Compile with –tcollect. | Automatically instruments all function entries/exits. | Requires recompile of code. |
| Add API calls to source code. | Can selectively instrument desired code sections. | Requires code modification. |

# Views and Charts

Helps navigating through the trace data and keep orientation

Every View can contain several Charts

All Charts in a View are linked to a single:

- time-span
- set of threads
- set of functions

All Charts follow changes to View (e.g. zooming)



**Chart**

# Event Timeline



Get detailed impression of program structure

Display functions, messages, and collective operations for each rank/thread along time-axis

Retrieval of detailed event information

# Quantitative Timeline

Get impression on parallelism and load balance

Show for every function how many threads/ranks are currently executing it

# Flat Function Profile

## Statistics about functions

# Call Tree and Call Graph

## Function statistics including calling hierarchy

- Call Tree shows call stack

- Call Graph shows calling dependencies

# Communication Profiles

Statistics about point-to-point or collective communication

Matrix supports grouping by attributes in each dimension

- Sender, Receiver, Data volume per msg, Tag, Communicator, Type

## Available attributes

- Count, Bytes transferred, Time, Transfer rate

# Zooming

# Grouping and Aggregation

Allow analysis on different levels of detail by aggregating data upon group-definitions

## Functions and threads can be grouped hierarchically

- Process Groups and Function Groups

 All_Processes     Major Function Groups

## Arbitrary nesting is supported

- Functions/threads on the same level as groups

- User can define his/her own groups

## Aggregation is part of View-definition

- All charts in a View adapt to requested grouping

- All charts support aggregation

# Aggregation Example

# Tagging and Filtering

Help concentrating on relevant parts

Avoid getting lost in huge amounts of trace data

Define a set of interesting data

- E.g. all occurrences of function x

- E.g. all messages with tag y on communicator z

Combine several filters:
Intersection, Union, Complement

Apply it

- Tagging: Highlight messages

- Filtering: Suppress all non-matching events

# Tagging Example

# Filtering Example

# MPI Performance Assistance

Automatic Performance Assistant

Detect common MPI performance issues

Automated tips on potential solutions



Automatically detect performance issues and their impact on runtime

# Summary page shows computation vs. communication breakdown

# MPI-3.0 Support

## Support for major MPI-3.0 features

- Non-blocking collectives

- Fast RMA

- Large counts



Non-blocking Allreduce (MPI_Iallreduce)

# Legal Disclaimer & Optimization Notice

**Optimization Notice**