



STATISTICS TOOL KIT

Инструментальные средства
исследования эффективности
параллельных
приложений

Докладчик:

Новаев Д.А.

Содокладчики:

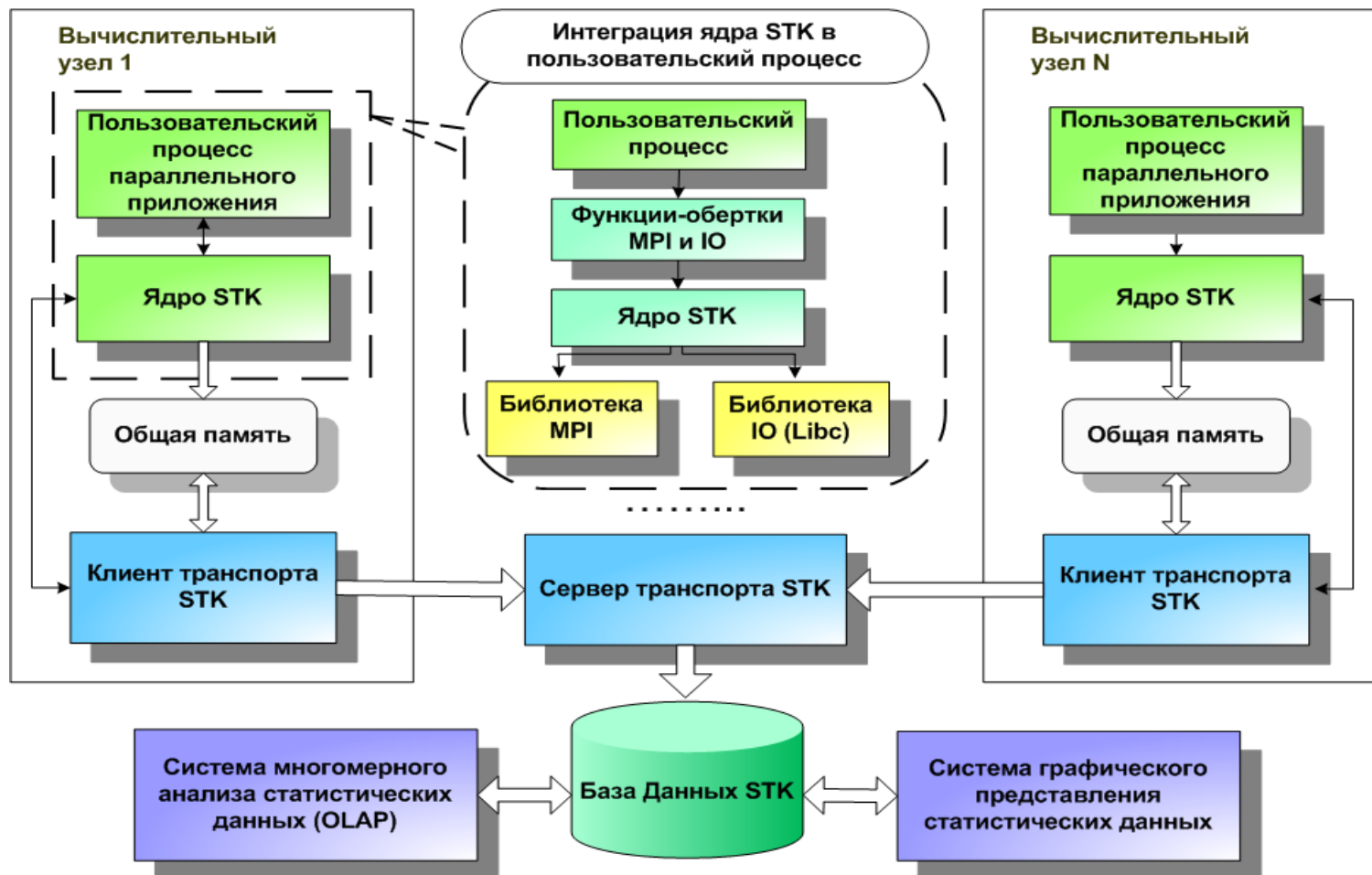
Колпаков С.И.,

Липов Д.И.

Назначение

- Сбор и представление данных по эффективности выполнения параллельного производственного счета на компонентах НВК ИТМФ ВНИИЭФ.
- Организация контроля за эффективностью выполнения параллельных задач пользователей и обеспечение более эффективного счета задач.
- Предоставление разработчикам параллельных программных комплексов средств для определения признаков и причин неэффективной работы как всей программы в целом, так и ее отдельных фрагментов.

Архитектура STK



Особенности STK



Проблемно-ориентированные библиотеки

Сторонние ИТМФ РФЯЦ-ВНИИЭФ

Название	Назначение	Поддержка в STK
<i>УРС-ОФ</i>	Единая унифицированная система расчета теплофизических свойств веществ	Есть
<i>PMLP/Parsol</i>	Параллельные решатели	Возможна
<i>MARX</i>	Расчет групповых анизотропных макроскопических констант среды	Возможна
<i>EFR</i>	Универсальное представление расчетных данных	Возможна
... и другие, написанные на C/C++ или Fortran 77/90 для ОС Unix		
<i>MPI (Mvarich, OpenMPI, HPMPI и др.)</i>	Передача сообщений	Есть
<i>UNIX Libc</i>	Ввод-вывод	Есть
<i>CUDA</i>	Расчеты на графических ускорителях NVIDIA	Есть
... и другие, написанные на C/C++ или Fortran 77/90 для ОС Unix		

Группы перехватываемых вызовов

MPI	Операции для инициализации и завершения среды MPI
	Блокирующие операции передачи
	Блокирующие операции приема
	Неблокирующие операции передачи
	Неблокирующие операции приема
	Операции приема-передачи
	Операция барьерной синхронизации
	Операции для синхронизации процессов при обменах
	Коллективные операции обменов
	Возобновляемые (persistent) запросы на передачу
	Возобновляемые (persistent) запросы на прием
	Ю
Операции чтения из файлов	
Операции записи в файлы	
Операции со стандартным выводом	
Операции для работы с метаданными файлов (время создания, размер, права доступа, владелец и т.д.)	
УРС-ОФ	Подпрограммы инициализации
	Подпрограммы расчетов

Характеристики выполнения групп вызовов:

1. **Utime** - время ЦП выполнения группы вызовов в режиме пользователя
2. **Stime** - время ЦП выполнения группы вызовов в режиме ядра
3. **Count_Call** - количество вызовов в группе, перехваченных ядром STK
4. **Walltime** - календарное время выполнения группы вызовов
5. **Max_Walltime** - максимальное календарное время выполнения функции в группе вызовов
6. **Sizebuf** - объем данных для операций передачи, файловых операций
7. **Max_Sizebuf** - максимальный объем данных при передаче между процессами пользовательской задачи для группы MPI вызовов

Оценка эффективности выполнения параллельных приложений

$$E = \frac{\sum_{i=1}^N T_{Calc_i}}{\sum_{i=1}^N T_{All_i}} \cdot 100\% \quad (1), \text{ где}$$

E - показатель эффективности выполнения параллельного приложения

N - количество процессов распараллеливания

T_{Calc} - время арифметических вычислений

T_{All} - общее календарное время выполнения

$$T_{Calc} = T_{All} - T_{MPI} - T_{IO} - T_{Idle_Cpu} \quad (2), \text{ где}$$

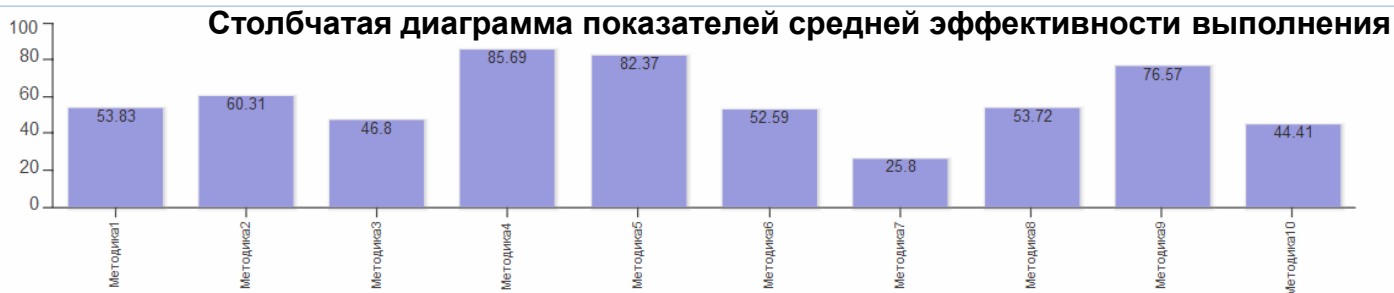
T_{MPI} - календарное время на MPI

T_{IO} - календарное время на ввод-вывод

T_{Idle_Cpu} - время простоя ЦП без учета MPI и ввода-вывода

Средства визуализации данных STK

Statistics Tool Kit (Статистика по эффективности выполнения параллельных приложений)



Интегральная статистика:

Фильтр ▾

Название	Статус	Исп.время, ч	Исп.время ЦП, ч	Ср.эфф-ть %	MPI %	IO %	URS-OF %
Кластер1		5076		68.63	27.51	1.13	0.01
Методика1		3685		53.83	46.06	0.1	10.15
Иванов И.И.		2655					
Процессов: 24 Тема: 04969.000 Заказчик: Сидоров С.С.		2290					
ID_задания: 20344 Завершена: 01.09.2011		2290					
Процессов: 60 Тема: 04969.000 Заказчик: Сидоров С.С.		365.27	21894.23				
ID_задания: 20544 Завершена:	В счете	365.27	21894.23				
Петров П.П.		1030.56	74126.38				
Процессов: 72 Тема: 04969.000		1030.56	74126.38				
Методика2		893.61	32137.63				
Методика3		388.69	18638.36				
Методика4		15.19	697.87				
Методика5		14.98	359.25				
Методика6		0.2	1.29				
Методика7		0.09	2.21				
Методика8		17.23	1101.76				
Методика9		36.81	2206.47				
Методика10		24.29	388.21				
Кластер2		2489.34	49720.23				

Комм. нагрузки MPI
Статистика счета
Графики
Диагностика взаимоблокировки
"Подозрительные" узлы и процессы

Иерархическое представление данных на 5 уровнях вложенности:

- 1) Многопроцессорные Вычислительные Системы;
- 2) Методики (Названия параллельных программных комплексов) или Группы пользователей;
- 3) Исполнители;
- 4) Количество процессов распараллеливания;
- 5) Задания исполнителя.

Использование библиотек:

Название: Кластер1

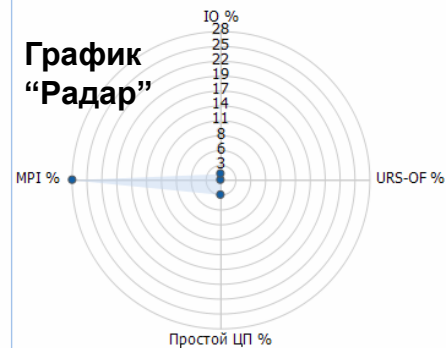
MPI %: 27.51 **Показатели (%)**

IO %: 1.13 **простая ЦП,**

URS-OF %: 0.01 **использования**

Доля простая ЦП %: 2.72 **библиотек**

Набор URS-OF констант: F074F092, F076F010, U005F012, U005F095, U005F097



Табличное представление данных

- Характеристики использования графических ускорителей (GPUs), библиотек MPI, IO и УРС-ОФ;
- Характеристики выполнения процессов параллельного приложения (использование памяти, время выполнения, показатели эффективности, низкоуровневые метрики PAPI и др.).

Статистика счета задания

Хар-ки использования процессами осн. ресурсов		Хар-ки использования MPI		Хар-ки использования IO		Хар-ки использования URS-OF		Хар-ки использования GPUs	
Номер пр.	Имя узла	Время ЦП пользо-ва	Время ЦП ядра, ми	Календарн	Объем занимаемой	Эффективность, %	MPI, %	IO, %	URS-OF, %
0	ve2524	Sort Ascending	0.08	6	3.16	57.35	42.04	0.7	0
1	ve2524	Sort Descending	0.1	6	3.16	75.13	24.36	0.72	0
2	ve2524		0.1	6	3.16	96.9	2.53	0.74	0
3	ve2524	Columns							
4	ve2524								
5	ve2524								
6	ve2524	5.82							
7	ve2524	5.81							
8	ve2524	5.81							
9	ve2524	5.79							
10	ve2524	5.79							
11	ve2524	5.82							
12	ve2525	5.78							
13	ve2525	5.83							
14	ve2525	5.78							
15	ve2525	5.79							
16	ve2525	5.78	0.11	5.99	3.16	48.43	51	0.77	0
17	ve2525	5.81	0.09	5.99	3.16	69.35	30.13	0.76	0
18	ve2525	5.77	0.11	5.99	3.16	47.69	51.67	0.76	0
19	ve2525	5.78	0.11	5.99	3.16	39.04	60.46	0.77	0

Page 1 of 6

Строки: [1 - 20] из 120

Табличное представление данных

- Коммуникационная нагрузка (темпы вызовов и объемы передач) со стороны MPI:

Коммуникационная нагрузка MPI

Интегральная статистика Детальная статистика

Интегральная коммуникационная нагрузка MPI

Номер процесса	Темп двухточечных обмен	Темп коллективных обмен	Темп Barrier, ед/с	Темп синхр.обменов, ед/с	Размер двухточечных обм	Размер коллективных обм	Темп всех MPI, ед/с	
Макс	79655.22	505.77	0	Sort Ascending	0	3228.37	34348.8	80160.99
Средн	1592.6	303.37	0.0	Sort Descending	0	1917.13	572.48	1895.96
Мин	780.95	291.13	0.0		0	718.47	0	1079.93
Сумма	127111.41	18291.36	3.4	Columns			34348.82	145402.8

Коммуникационная нагрузка MPI

Интегральная статистика Детальная статистика

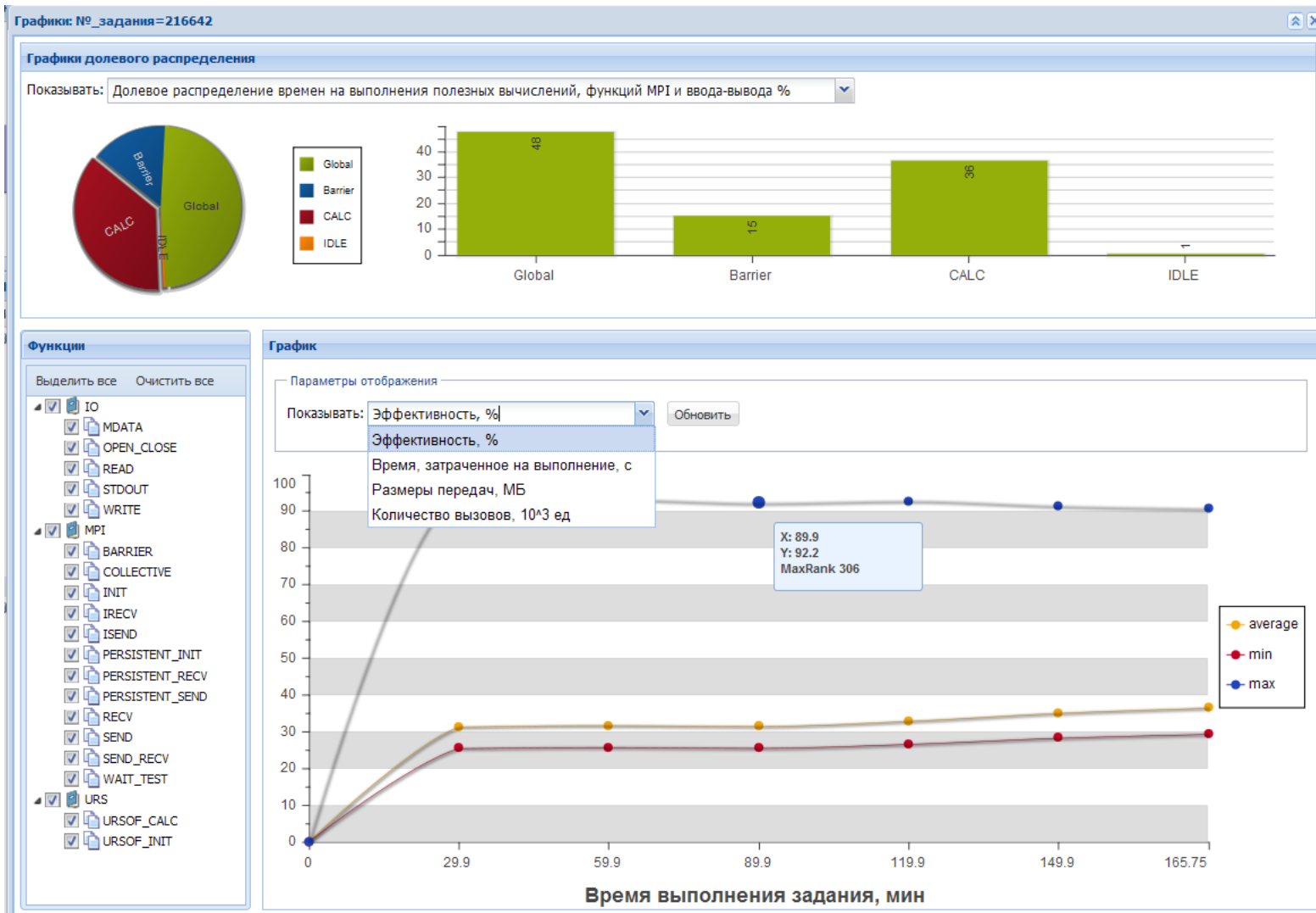
Коммуникационная нагрузка MPI по каждому процессу

Номер процесса	Темп двухточечных обмен	Темп коллективных обмен	Темп Barrier, ед/с	Темп синхр.обменов, ед/с	Размер двухточечных обм	Размер коллективных обм	Темп всех MPI, ед/с
0	79655.22	505.77	0.1	0	718.47	34348.8	80160.99
1	810.56	291.54	0.06	0	3038.87	0	1102.1
2	809.64	295.03	0.06	0	3063.76	0	1104.67
3	834.64	300.53	0.06	0	3040.49	0	1135.17
4	826.73	297.11	0.06	0	3036.6	0	1123.85
5	797.8	291.13	0.06	0	3065.42	0	1088.93
6	817.41	300.08	0.06	0	3077.16	0	1117.49
7	810.86	293.22	0.06	0	3047.73	0	1104.07
8	830.51	300.56	0.06	0	3049.34	0	1131.07
9	804.79	293.37	0.06	0	3063.25	0	1098.16
10	796.59	293.86	0.06	0	3086.01	0	1090.46
11	811.23	295.04	0.06	0	3058.61	0	1106.27
12	805.01	295.51	0.06	0	3076.27	0	1100.52
13	811.03	295.55	0.06	0	3062.63	0	1106.58
14	818.8	299.16	0.06	0	3067.57	0	1117.96
15	805.85	299.2	0.06	0	3098.96	0	1105.05
16	802.3	296.51	0.06	0	3090.06	0	1098.81
17	808.99	296.69	0.06	0	3074.76	0	1105.68
18	798.45	292.87	0.06	0	3075.28	0	1091.32
19	817.61	299.61	0.06	0	3073.39	0	1117.22

Page 1 of 3 Строки: [1 - 20] из 60

- Номер процесса
- Темп двухточечных обмен, ед/с
- Темп коллективных обмен, ед/с
- Темп Barrier, ед/с
- Темп синхр.обменов, ед/с
- Размер двухточечных обмен, Б
- Размер коллективных обмен, Б
- Темп всех MPI, ед/с
- Темп Bcast, ед/с
- Темп Scatter, ед/с
- Темп Gather, ед/с
- Темп AllGather, ед/с
- Темп All2All, ед/с
- Темп AllReduce, ед/с
- Темп Reduce, ед/с
- Темп Reduce-Scatter, ед/с
- Размер Bcast, Б
- Размер Scatter, Б
- Размер Gather, Б
- Размер AllGather, Б
- Размер All2All, Б
- Размер AllReduce, Б
- Размер Reduce, Б
- Размер Reduce-Scatter, Б

Графическое представление данных



Параметры использования GPU

1. **GPU_Usage** - процент времени относительно прошедшей секунды, в течение которого на графическом ускорителе выполнялась одна или несколько функций для GPU (kernels);
2. **Max_GPU_Usage** - пиковая загрузка графического ускорителя, полученная в течение времени выполнения параллельного приложения;
3. **Mem_Usage** - процент времени относительно прошедшей секунды, в течение которого глобальная (device) память GPU была прочитана или записана;
4. **Max_Mem_Usage** - пиковая загрузка при работе (чтение/запись) с памятью графического ускорителя, полученная в течение времени выполнения параллельного приложения;
5. **Mem_Size** - суммарное количество памяти, зарезервированное всеми активными каналами на GPU;
6. **Max_Mem_Size** - пиковое значение объема используемой памяти графического ускорителя, полученное в течение времени выполнения параллельного приложения.

Аппаратные счетчики современных микропроцессоров

- Набор регистров, хранящих число возникших событий:
 - *Общее число тактов;*
 - *Общее число команд;*
 - *Число операций с плавающей точкой;*
 - *Число тактов простоя функционального устройства;*
 - *Число промахов при работе с кэш-памятью;*
 - *Число промахов при работе с виртуальной памятью.*
- Отслеживание этих событий облегчает оптимизацию выполнения программы на данном CPU.

Интегральные параметры по событиям аппаратных счетчиков, представляемые STK

1. ***Мфлопс***, ед/сек - число миллионов операций с плавающей точкой в секунду;
2. ***Флоп***, ед. - число операций с плавающей точкой;
3. ***Обращение к кэшу L2***, ед. - число обращений в кэш второго уровня;
4. ***Промахов кэша L3***, ед. - число промахов кэша третьего уровня;
5. ***К эфф.памяти*** (число обращений к кэшу 2-го уровня/ число промахов кэша 3-го уровня) - коэффициент эффективности использования памяти, показывающий сколько обращений к кэшу второго уровня приходится на 1 промах кэша третьего уровня.

Заключение

- ✓ STK установлена на всех основных компонентах ВЦ РФЯЦ-ВНИИЭФ;
- ✓ STK входит в БСППО ВНИИЭФ и установлен на КС-ЭВМ разработки ВНИИЭФ, поставляемых в ряд ведущих предприятий авиа-, авто-, космической, атомной энергетики России;
- ✓ Перспективы развития:
 - сбор и представление данных по эффективности выполнения большого количества процессов (до 1000 000) параллельных задач;
 - 3D-визуализация VSTK-3D (online);
 - Интегральная статистика СУЗ (оценки работы/простоя выч.систем, параметры заданий и др.)

Спасибо за внимание!