

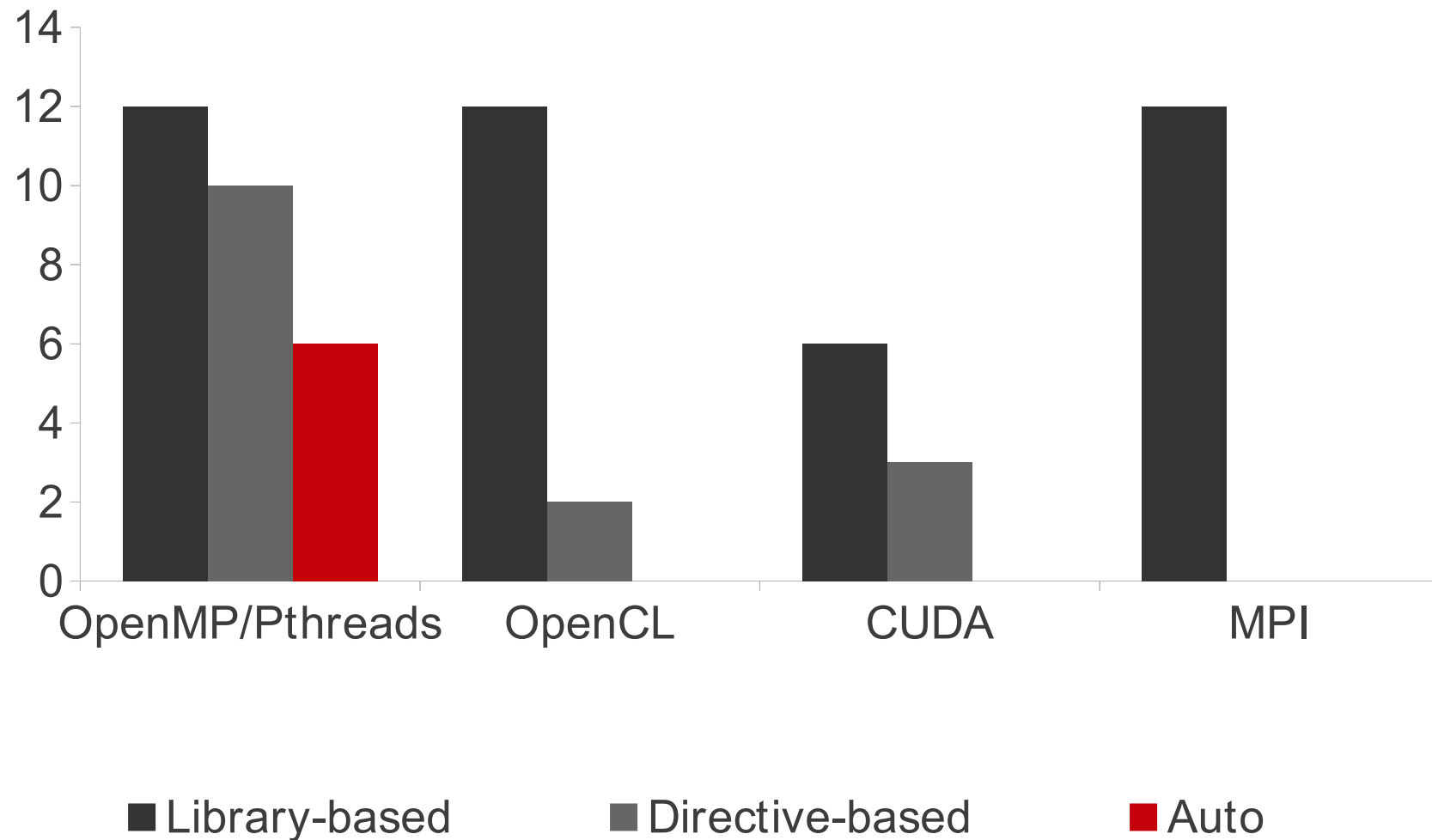
# Practical Efficiency of Optimizing Compilers in Parallel Scientific Applications

Bogdanov A.V., Gankevich I.G.  
Saint-Petersburg State University  
[bogdanov@csa.ru](mailto:bogdanov@csa.ru), [gig.spb@gmail.com](mailto:gig.spb@gmail.com)

# Outline

- Compilers overview
- Micro benchmarks results
- Parallel benchmarks results
- Conclusion

# Parallel standards support



Compilers being analyzed: GCC, Intel, PGI, Clang, Lahey, NAG, CAPS HMPP, Pathscale, Absoft, Silverfrost FTN95, MSVC, IBM XL

# Compilers being benchmarked

- GCC, Intel, PGI
  - All support OpenMP
  - All support autoparallelization on CPU
  - PGI supports directive based parallelization on GPU

# Parallelization benchmark setup

- Source code
  - Reference BLAS implementation ([www.netlib.org](http://www.netlib.org)) compiled with automatic parallelization options
  - Hand-tuned Goto BLAS implementation compiled with library-chosen optimization options
- Target platform
  - HP SL390s G7
  - 2x Intel X5650 2.67 Ghz (12 cores)
  - 96 Gb RAM
  - 3x NVIDIA Tesla M2050

# Parallelization benchmark setup

- Routines being analyzed
  - SGEMM, SSYMM, SSYR2K, SSYRK, STRMM, STRSM
- Variations
  - Different sizes (512, 1024, 2048, 4096, 8192)
  - Transposed/nontransposed matrices
  - Upper/lower triangular matrices
  - Left/right sided
- A total of 128 test cases were analysed for each compiler and different amount of threads

# Average performance

	Hand-tuned Goto BLAS			Autoparallelized BLAS (CPU)			Directive-parallelized BLAS (GPU)
	Intel	PGI	GCC	Intel	PGI	GCC	PGI
SGEMM	111974	<b>130197</b>	57816	<b>16173</b>	7333	1824	1939
SSYMM	108611	<b>130466</b>	60525	9089	<b>12123</b>	2390	1632
SSYR2K	103813	<b>112828</b>	60230	<b>18157</b>	13728	3136	1917
SSYRK	95721	<b>106000</b>	43239	<b>14562</b>	10027	1811	1218
STRMM	99073	<b>115354</b>	43204	<b>14423</b>	7148	1990	2102
STRSM	97371	<b>115009</b>	41857	<b>11381</b>	6295	1745	1998

Performance is reported in Mflops as the average of all routine arguments variations

# SGEMM source code

SGEMM( $A^T$ , B)	SGEMM( $A^T$ , $B^T$ )
<pre>DO 120 J = 1,N   DO 110 I = 1,M     TEMP = ZERO     DO 100 L = 1,K       TEMP = TEMP + A(L, I) * B(L, J)     CONTINUE     IF (BETA.EQ.ZERO) THEN       C(I, J) = ALPHA*TEMP     ELSE       C(I, J) = ALPHA*TEMP + BETA*C(I, J)     END IF     CONTINUE CONTINUE</pre>	<pre>DO 200 J = 1,N   DO 190 I = 1,M     TEMP = ZERO     DO 180 L = 1,K       TEMP = TEMP + A(L, I) * B(J, L)     CONTINUE     IF (BETA.EQ.ZERO) THEN       C(I, J) = ALPHA*TEMP     ELSE       C(I, J) = ALPHA*TEMP + BETA*C(I, J)     END IF     CONTINUE CONTINUE</pre>



# Memory performance

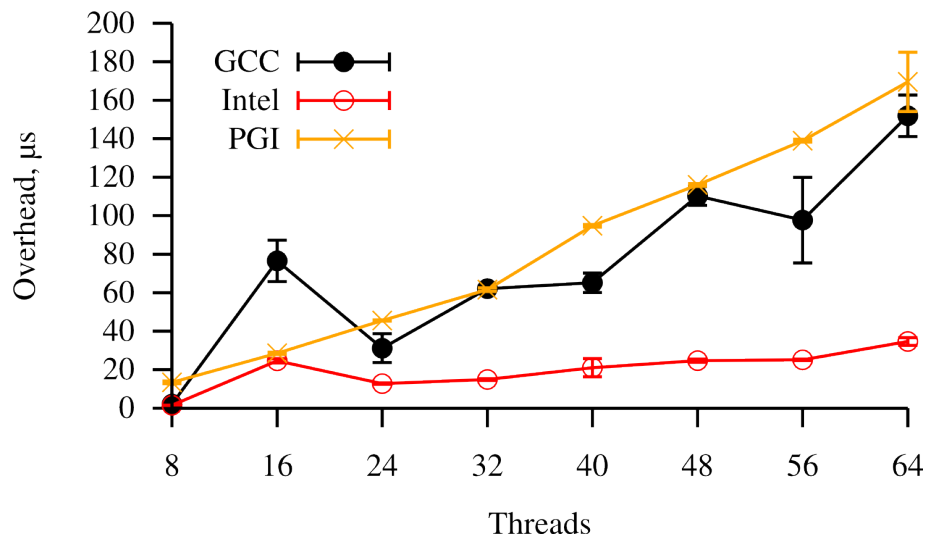
	Reads $\times 10^9$	Writes $\times 10^9$	Instructions $\times 10^9$	Misses $\times 10^9$
SGEMM( $A^T$ , B)	137.4	0.016	412	0.008
SGEMM( $A^T$ , $B^T$ )			549	0.104

# OpenMP benchmark setup

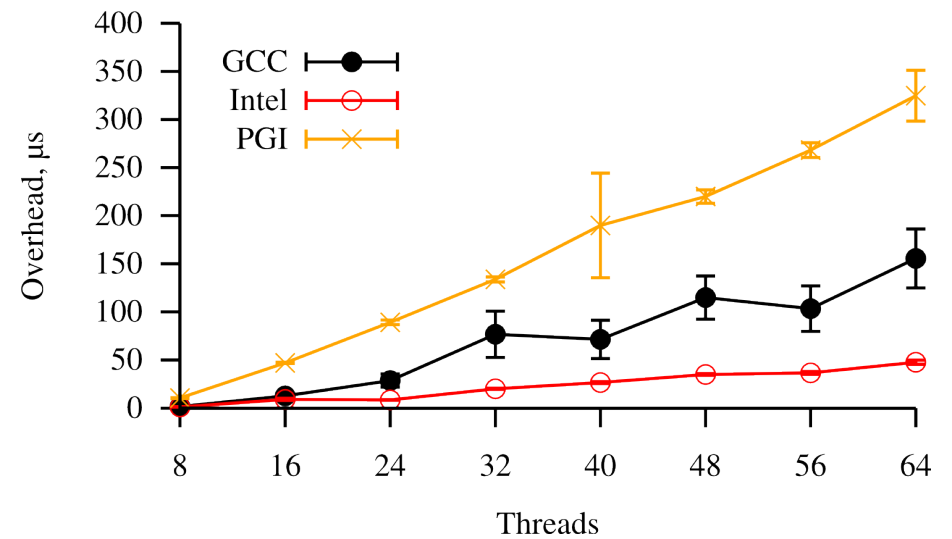
- Source code
  - EPCC benchmarks compiled with maximum optimization level
- Target platform
  - HP Proliant DL980
  - 8x Intel X7560 2.2 Ghz (64 cores)
  - 512 Gb RAM

# OpenMP directives overhead

#pragma omp reduction



#pragma omp for schedule(dynamic)



# Conclusions

- Autoparallelization on CPU is suitable for prototyping
- Directive-based parallelization on GPU is immature and requires hand tuning

Thank you