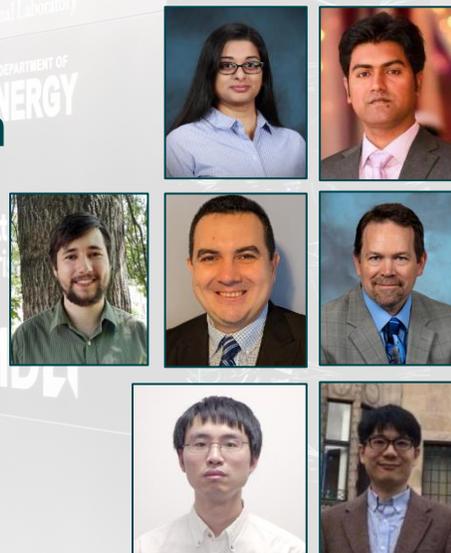


ChatHPC: Building the foundation of an AI assisted and productive HPC ecosystem



PRESENTED BY

Pedro Valero-Lara

PhD. Sr Computer Scientist <valerolarap@ornl.gov>



ORNL IS MANAGED BY UT-BATTELLE LLC
FOR THE US DEPARTMENT OF ENERGY

FRONTIER



ChatHPC, What is that?

Infrastructure:

- A CLI Python library on top of:
 - Code Llama, PyTorch, and LoRA
- The ChatHPC library makes:
 - Fine-tuning, Testing, Refinement, Merge, and Inference

Ecosystem:

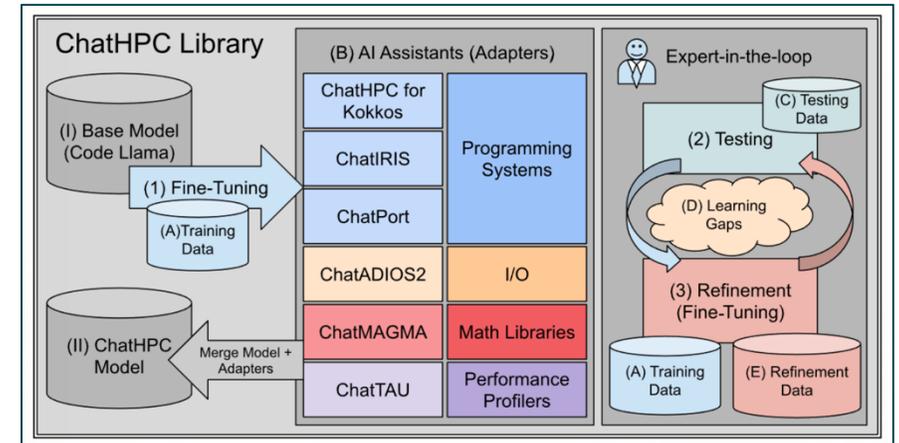
- A collection of AI assistants leveraging existing HPC efforts for multiple domains:
 - Programming Models, Math Libraries, I/O, Tooling, etc.
- And priorities:
 - Parallelization, Portability, Optimization, Evaluation, etc.

ChatHPC CLI:

```
1 $ chathpc train # Finetune the assistant.
2 $ chathpc verify # Verify the assistant on training set.
3 $ chathpc test # Test the assistant on unseen data.
4 $ chathpc run # Interactively run the assistant.
```

Interactive Run Session:

```
1 $ chathpc ()> /context
2 Context: Introduction to Kokkos
3 $ chathpc (Introduction to Kokkos)> What is LayoutLeft?
4 LayoutLeft refers column-major layout where consecutive entries in
the same column of a 2-D array are contiguous in memory.
```

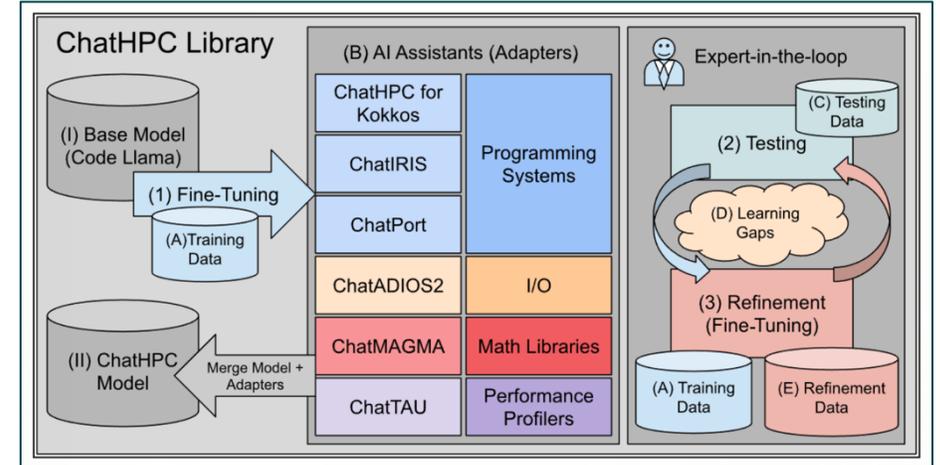


Design and Implementation

A simple 3-step iterative process:

1. Fine-tuning -> AI assistant (optimizers)
 - Training Data
 - The better the data the better the model
 - It is more about quality than quantity
2. Testing -> Learning gaps
 - Testing Data != Training Data
 - Testing Data may contain question about non-fine-tuned capabilities
3. Refinement -> More capable AI assistants
 - Expert-Supervised Refinement Data

AI assistants can be merge into a single LLM: ChatHPC model



```
1 [
2 {
3   "question": "Can you translate this HIP code to Kokkos?
4   __global__ void stencil (int n, float dx, float dy, float *u,
5     float *u_new, float *f) {
6     int i = blockIdx.x * blockDim.x + threadIdx.x;
7     int j = blockIdx.y * blockDim.y + threadIdx.y;
8     if(i > 0 && i < n-1 && j > 0 && j < n-1){
9       u_new[i][j] = 0.25 * (u[i-1][j] + u[i][j+1]
10        + u[i][j-1] + u[i+1][j]) + f[i][j];
11     }
12   }",
13   "context": "HIP to Kokkos translation",
14   "answer":
15     "Kokkos::parallel_for (\" stencil\",
16       Kokkos::MDRangePolicy<Kokkos::Rank<2>>
17         ({1, 1}, {n-1, n-1}),
18       KOKKOS_LAMBDA(const int i, const int j) {
19         u_new(i, j) = 0.25 * (u(i-1, j) + u(i, j+1)
20           + u(i, j-1) + u(i+1, j) + f(i, j));
21       });"
22 }
```

Analysis

- Code Llama 7B [13GB]
 - ChatGPT 4o 1.8T -> 257x less parameters
- AI assistants (~100MB)
- Training data (~KB)
- Evaluation data (~KB)
- Fine-tuning:
 - No more than 15 min. in a one NVIDIA 2x H100 node
- Evaluation
 - % of prompts (evaluation data) passed
 - Evaluation data is different to training data
- Compare against
 - 7 billion parameter Meta's Code Llama base model (no fine-tuning)
 - 1.8 trillion parameter OpenAI's ChatGPT 4o model
- Data
 - Synthetic data in most of the AI-assistants created (supervised learning)
 - Kernels
 - Applications
 - Scripts
 - Applications:
 - Dense/Sparse Linear Algebra, PDEs, CFD, AI, Image Processing, ...
 - Expert-in-the-loop
 - If we want AI responding as expert, we need data provided by experts
 - Create data
 - Identify learning gaps
 - Minimize data demanding
 - Self-learning
 - Reinforcements, if necessary (learning gaps)
 - Which targets are more demanding?
 - Not use all data (training data) we can use

ChatHPC for Kokkos

Model/Context	CodeLlama	ChatHPC for Kokkos Initial	ChatHPC for Kokkos Refinement	ChatGPT
Documentation	9.5%	89.0%	-	81.0%
Installation	27.2%	78.0%	-	45.5%
Development	0.0%	85.0%	-	15.1%
Parallelization	0.0%	45.5%	90.9%	66.7%
Portability	0.0%	55.6%	85.9%	33.4%
CUDA*	0.0%	53.2%	87.2%	41.3%
OpenACC*	0.0%	58.7%	83.5%	31.2%

Kokkos: A C++ API portable library for parallel programming

Parallelization and Portability

- Training (HPC and scientific kernels)
 - BLAS, Sparse, Euler methods, Conjugate Gradient, LBM, AI convolutions
- Learning gaps
 - Sparse Computation and Iterative Solvers
 - 12% more data

No CUDA and OpenACC codes used for fine-tuning

Translate the next CUDA code to Kokkos:

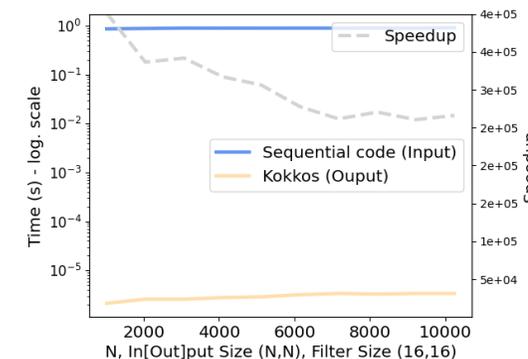
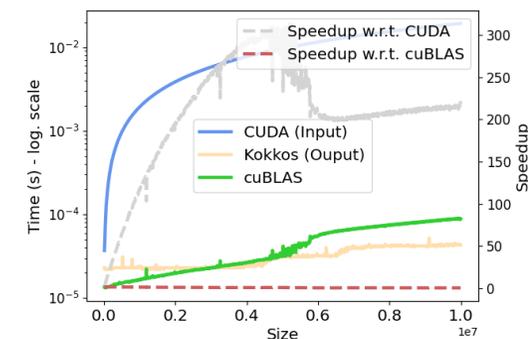
```

1 __global__ void sdot(int n, float *x, float *y, float *result) {
2   int i = threadIdx.x + blockIdx.x * blockDim.x;
3   float update = 0.0;
4   if (i < SIZE) {
5     update = x[i] * y[i];
6     atomicAdd(result, update);
7   }
8 }
    
```

ChatKokkos output:

```

1 Kokkos::parallel_reduce(SIZE,
2   KOKKOS_LAMBDA(int i, float &update){
3     update += x(m) * y(m);
4   },
5   Kokkos::Sum<float>(result));
    
```



ChatADIOS2

Model/Context	CodeLlama	ChatADIOS2 Intial	ChatADIOS2 Refinement	ChatGPT
Introduction	0.0%	77.0%	100.0%	33.3%
Variable Definition	0.0%	50.0%	100.0%	33.3%
Data Compression	0.0%	25.0%	100.0%	0.0%
Parallelization	0.0%	0.0%	100.0%	100.0%

ADIOS2: A portable and HPC I/O library

Capabilities:

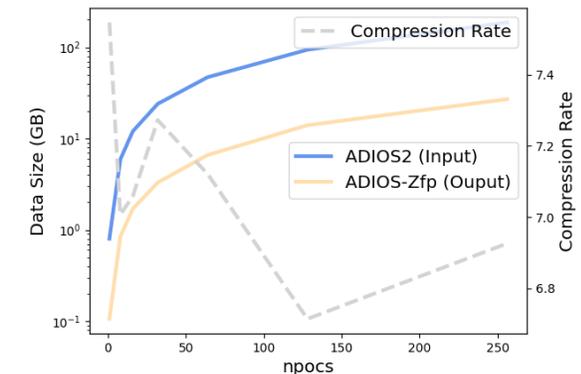
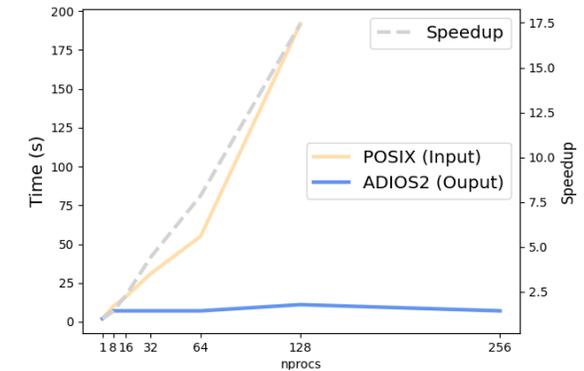
- Q&A on API, Variable definition, Compression, and Parallelization (serial POSIX I/O to ADIOS2)

Learning gaps:

- Wrong use of ADIOS2 API
 - Local and global N-dimension variables
 - Compression factor (parameters)

```

1 adios2::ADIOS adios("adios2_config.yaml", mpiComm);
2 adios2::IO io = adios.DeclareIO("WriteIO");
3 adios2::Variable<double> varU =
4     io.DefineVariable<double>("U", shape, start, count);
5 adios2::Operator Zfp =
6     adios.DefineOperator("Zfp", adios2::ops::LossyZFP);
7 varU.AddOperation(Zfp, {{ adios2::ops::zfp::key::rate, "9" }});
8
9 adios2::Engine bpWriter =
10    io.Open("output.bp", adios2::Mode::Write);
11 while( run_simulation ){
12     bpWriter.BeginStep();
13     bpWriter.Put<double>(varU, U.data());
14     ...
15     bpWriter.EndStep();
16 }
17 bpWriter.Close();
    
```



ChatIRIS

Model/Context	CodeLlama	ChatIRIS Intial	ChatGPT
Documentation	0.0%	95.0%	70.0%
OpenMP -> IRIS	0.0%	82.5%	10.0%
CUDA -> IRIS	0.0%	90.0%	50.0%
HIP -> IRIS*	0.0%	81.6%	10.0%

IRIS: A Task-based programming model for heterogeneous computing

- Taskification of applications
 - Deploy applications on any architecture using multiple devices
- Portability
 - OpenMP, CUDA, HIP to IRIS task code
- Training
 - Basic applications (AXPY)
- Testing
 - More complex applications with multiple kernel invocations
- No refinements were needed
- HIP was not used for fine-tuning

Translate the next CUDA code to IRIS:

```
1 alpha = 2.0;
2 X = (float *) malloc(SIZE * sizeof(float));
3 Y = (float *) malloc(SIZE * sizeof(float));
4 Z = (float *) malloc(SIZE * sizeof(float));
5 for (int i = 0; i < SIZE; i++) { X[i] = i; Y[i] = i;}
6 cudaMalloc((void*)&d_x, SIZE * sizeof(float));
7 cudaMalloc((void*)&d_y, SIZE * sizeof(float));
8 cudaMalloc((void*)&d_z, SIZE * sizeof(float));
9 // Initialization of arrays
10 cudaMemcpy(d_x, x, SIZE * sizeof(float), cudaMemcpyHostToDevice);
11 cudaMemcpy(d_y, y, SIZE * sizeof(float), cudaMemcpyHostToDevice);
12 int threads = 512;
13 int blocks = floor(SIZE/threads) + 1;
14 saxpy<<<blocks, threads>>>(alpha, d_x, d_y, SIZE);
15 cudaDeviceSynchronize();
16 sdot<<<blocks, threads>>>(d_z, d_x, SIZE);
17 cudaDeviceSynchronize();
18 cudaMemcpy(z, d_z, SIZE * sizeof(float), cudaMemcpyDeviceToHost);
19 cudaFree(d_x);
20 cudaFree(d_y);
21 cudaFree(d_z);
22 free(X);
23 free(Y);
24 free(Z);
```

ChatIRIS output:

```
1 // Architecture-agnostic IRIS task code
2 alpha = 2.0;
3 iris_init(&argc, &argv, 1);
4 iris_mem iris_x;
5 x = (float *) malloc(SIZE * sizeof(float));
6 iris_data_mem_create(&iris_x, x, SIZE * sizeof(float));
7 iris_mem iris_y;
8 y = (float *) malloc(SIZE * sizeof(float));
9 iris_data_mem_create(&mem_Y, Y, SIZE * sizeof(float));
10 iris_mem iris_z;
11 z = (float *) malloc(SIZE * sizeof(float));
12 iris_data_mem_create(&mem_z, z, SIZE * sizeof(float));
13 // Initialization of arrays
14 iris_graph graph;
15 iris_graph_create(&graph);
16 iris_task saxpy_task;
17 iris_task sdot_task;
18 iris_task_create(&saxpy_task);
19 iris_task_create(&sdot_task);
20 void* saxpy_params[3]={&alpha, &iris_x, &iris_y};
21 int saxpy_params_info[3]={sizeof(A), iris_w, iris_r};
22 iris_task_kernel(task0, "saxpy", 1, NULL, &SIZE, NULL, 3,
23 saxpy_params, saxpy_params_info);
24 void* sdot_params[2]={&iris_z, &iris_x};
25 int sdot_params_info[2]={iris_w, iris_r};
26 iris_task_kernel(task1, "sdot", 1, NULL, &SIZE, NULL, 2,
27 sdot_params, sdot_params_info);
28 iris_graph_task(graph, saxpy_task, iris_default, NULL);
29 iris_graph_task(graph, sdot_task, iris_default, NULL);
30 iris_graph_submit(graph, iris_any, 1);
31 free(X);
32 free(Y);
33 free(Z);
34 iris_finalize();
```

ChatMAGMA



Model/Context	CodeLlama	ChatMAGMA Initial	ChatMAGMA Refinement	ChatGPT
MKL -> MAGMA	0.0%	37.5%	93.0%	0.0%
cuBLAS/Solver -> MAGMA	0.0%	50.0%	96.5%	0.0%
hipBLAS/Solver -> MAGMA	0.0%	45.3%	95.3%	0.0%

MAGMA: A portable (CPU+GPU) BLAS and LAPACK library with +400 routines

Capabilities:

- Porting of vendor (arch)-specific applications to MAGMA
 - Support for heterogeneous (CPU+ GPU) computing
- Training
 - Well known BLAS (dgemm, dtrsm) and LAPACK (dgetsv, dpotrf) for double precision
- Testing
 - Rest of BLAS and LAPACK for other precisions
- Learning gaps:
 - MAGMA-specific data types (queues) for different precisions
 - 51.72% extra data

No AMD math libraries codes were used for fine-tuning

Translate the next Intel MKL code to MAGMA:

```
1 // Intel MKL (CPU-only) dgetrf code
2 double *A=(double *)mkl_malloc(n*n*sizeof(double),64);
3 int *ipiv=(int *)mkl_malloc(n*sizeof(int),64);
4 LAPACK_dgetrf(LAPACK_ROW_MAJOR,n,n,A,n,ipiv);
5 mkl_free(A);
6 mkl_free(ipiv);
```

Translate the next NVIDIA cuSolver code to MAGMA:

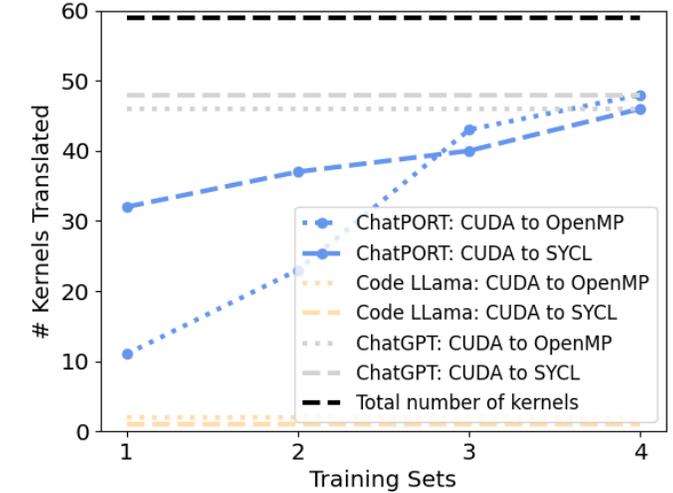
```
1 // NVIDIA cuSolver (GPU-only) dgetrf code
2 int lwork=0;
3 double *dwork;
4 int *dipiv;
5 double *dinfo;
6 cusolverDnHandle_t handle;
7 cusolverDnCreate(&handle);
8 double *hA=(double *)malloc(n*n*sizeof(double));
9 cusolverDnDgetrf_bufferSize(handle,n,n,dA,n,&lwork);
0 cudaMalloc((void **)&dA,n*n*sizeof(double));
1 cudaMalloc((void **)&dwork,sizeof(double)*lwork);
2 cudaMalloc((void **)&dipiv,n*sizeof(int));
3 cudaMalloc((void **)&dinfo,sizeof(int));
4 cudaMemcpy(dA,hA,n*n*sizeof(double),cudaMemcpyHostToDevice);
5 cusolverDnDgetrf(handle,n,dA,n,dwork,dipiv,dinfo);
6 cudaMemcpy(hA,dA,n*n*sizeof(double),cudaMemcpyDeviceToHost);
7 cudaFree(dA);
8 cudaFree(dwork);
9 cudaFree(dipiv);
0 cudaFree(dinfo);
1 free(hA);
2 cusolverDnDestroy(handle);
```

ChatMAGMA output:

```
1 // MAGMA (CPU+GPU) dgetrf code
2 magma_queue_t queue;
3 magma_device_t cdev;
4 magma_getdevice(&cdev);
5 magma_queue_create(cdev,&queue);
6 double *hA;
7 magmaDouble_ptr dA;
8 magma_int_t *ipiv;
9 magma_int_t info;
0 magma_dmalloc_cpu(&hA,n*n);
1 magma_dmalloc(&dA,n*n);
2 magma_imalloc_cpu(&ipiv,n);
3 magma_dsetmatrix(n,n,hA,n,dA,n,queue);
4 magma_dgetrf_gpu(n,n,dA,n,ipiv,&info);
5 magma_dgetmatrix(n,n,dA,n,hA,n,queue);
6 magma_free(dA);
7 magma_free_cpu(hA);
8 magma_free_cpu(ipiv);
9 magma_queue_destroy(queue);
```

ChatPORT(ing)

Training Set	#Kernels	Kernels
1	5	accuracy, geodesic, lr, maxpool3d, perplexity
2	10	damage, knn, heat2d, laplace3d, md
3	15	advCubatureHex3D, backprop, channelShuffle, lr, meanshift
4	22	chemv, clenergy, pathfinder, pointwise, swish, stddev, tissue



- Create an AI assistant for code portability from vendor-specific models (e.g., CUDA) to SYCL and OpenMP-offloading
- Data is provided by HeCBench
 - Training:
 - 4 different data sets with (5, 10, 15, and 22) kernels from multiple different domains
 - Testing
 - 59 CUDA kernels to be ported to OpenMP-offloading and SYCL

Remarks

Assistant	Code Llama	Initial	Extra Data	Refinement	ChatGPT
ChatHPC for Kokkos	7.34%	70.51%	12.33%	85.75	48.34%
ChatADIOS2	0.00%	40.87%	37.50%	100.00%	25.0%
ChatMAGMA	0.00%	44.26%	51.72%	94.93%	0.00%
ChatIRIS	0.00%	87.29%	-	-	35.00%
ChatPort	2.54%	36.44%	77.27%	79.66%	79.66%

- New HPC-specific capabilities on top of Code Llama across different HPC domains
 - Programming models, I/O, math libraries, tooling, and runtimes

And tasks:

- Parallelization, portability, performance analysis, scalability, etc.
- Provide higher levels (of up to 90% higher) of trustworthiness than ChatGPT quickly with modest resources
 - Increase data-size (expert-in-the-loop) during refinement learning has a higher impact on correctness
- More data (prompts) for testing than for training data
 - Testing data contains prompts about non-fine-tuned capabilities

Future Work

- **ChatMPI**

- A ChatHPC AI-assistant for MPI parallelization
 - 4x speedup w.r.t. MPI code generated by ChatGPT-5

- **3 Pillars**

- **Multimodality**

- Simplify the interaction with AI for HPC domains

- **Reasoning**

- Generate highly optimized HPC codes across domains and technologies

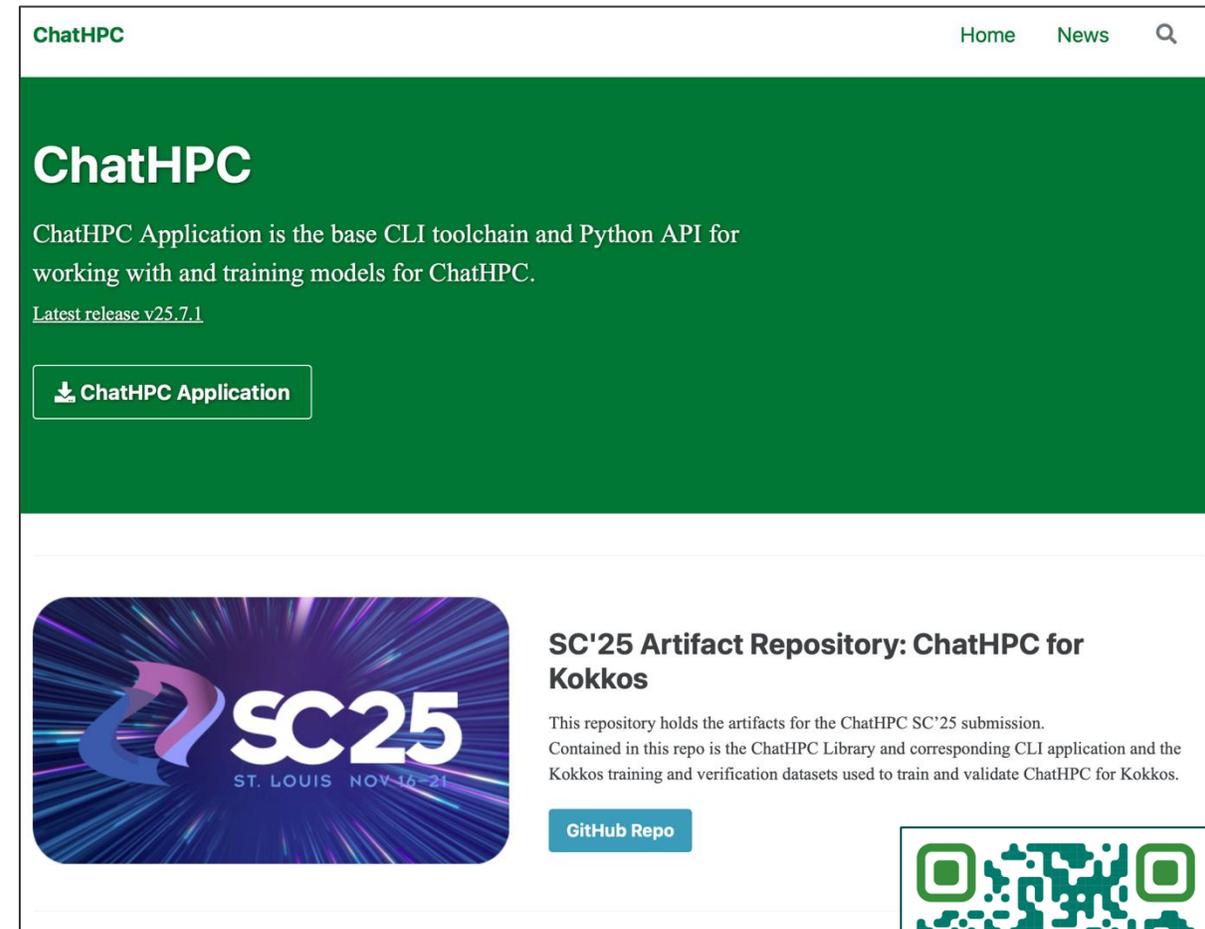
- **Agentic**

- Connect HPC and AI for complex problems



Resources and Community

- ChatHPC project:
<https://ornl.github.io/ChatHPC/>
- Meetings:
 - We organize a community meeting once a month.
 - Join us!
- Workshops:
 - LLM4HPC (co-located with ISC-HPC):
<https://ornl.github.io/events/llm4hpc2025/>
 - LLM4HPCAsia (collocated with SCA/HPC Asia)
<https://ornl.github.io/events/llm4hpcasia2026/>



ChatHPC

Home News

ChatHPC

ChatHPC Application is the base CLI toolchain and Python API for working with and training models for ChatHPC.

Latest release [v25.7.1](#)

↓ ChatHPC Application

SC'25 Artifact Repository: ChatHPC for Kokkos

This repository holds the artifacts for the ChatHPC SC'25 submission. Contained in this repo is the ChatHPC Library and corresponding CLI application and the Kokkos training and verification datasets used to train and validate ChatHPC for Kokkos.

GitHub Repo



Observations (Things to think about)

- **Fine-tuning**
 - Very effective, if you have the data
- **Trustworthiness**
 - You can elevate trustworthiness (BTW: better than ChatGPT)
- **Overfitting**
 - Is that a problem for domain-specific AI assistants?
- **Self-learning**
 - Train on simple problems and let LLMs to deal with more complex problems
 - Take advantage from similarities (similar patterns)
- **Synthetic or real data?**
 - Why not both?
- **Expert-in-the-loop**
 - Recommendable for domain-specific LLMs
- **Accessibility**
 - One NVIDIA GPU (better if you have 2), data (~KB), and 15 minutes
- **Simplify prompting**
 - No prompt engineering, more productive, non-experts can generate expert level responses
- **Impact on HPC**
 - Created capabilities for programming models, I/O, math libraries, performance tools, ...
 - **Help our community to use the HPC tools we have today**
- **Leveraging HPC efforts**
 - ChatHPC benefits from the big effort made during last years

Use ChatHPC!



ChatHPC: Building the foundation of an AI assisted and productive HPC ecosystem

Thanks!

PRESENTED BY

Pedro Valero-Lara

PhD. Sr Computer Scientist <valerolarap@ornl.gov>



ORNL IS MANAGED BY UT-BATTELLE LLC
FOR THE US DEPARTMENT OF ENERGY

Acknowledgments:

This research used resources of the Oak Ridge Leadership Computing Facility and the Experimental Computing Laboratory at the Oak Ridge National Laboratory, which is supported by the Office of Science of the US Department of Energy under Contract No. DE-AC05-00OR22725. This material is based upon work supported by the US Department of Energy, Office of Science's Advanced Scientific Computing Research program as part of the Advancements in Artificial Intelligence for Science program, Ellora and Durban projects, the Next Generation of Scientific Software Technologies program, S4PST and PESO projects, and the Scientific Discovery through Advanced Computing (SciDAC) program, RAPIDS2 SciDAC Institute for Computer Science, Data, and Artificial Intelligence.