ParEval-Repo: A Benchmark Suite for Evaluating LLMs with Repository-level HPC Translation Tasks

Josh H. Davis, Daniel Nichols, Ishan Khillan, Abhinav Bhatele







The Team





Joshua H. Davis



Daniel Nichols



Ishan Khillan



Abhinav Bhatele



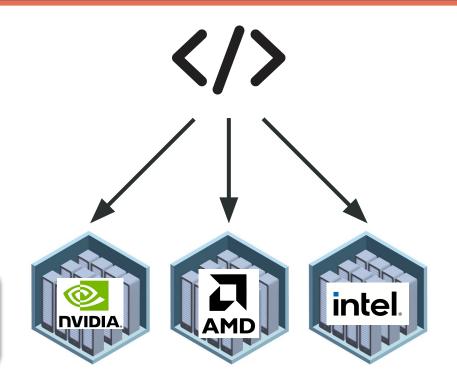


Why Automate Translation for HPC?



- Three different GPU vendors in top ten supercomputers
- Maximizing access requires using portable programming models
- But, manually porting is time-consuming and tedious!

Research Question: Can LLMs help to translate HPC code repositories into portable GPU programming models?



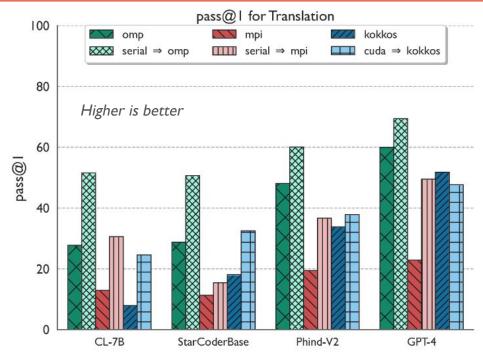




Background for LLM Parallel Code Translation



- ParEval found that parallel code tasks are uniquely difficult
- Function-level LLM translation of parallel code is more feasible than generation
- But we don't know if this will scale to full repositories



D. Nichols, J. H Davis, et al. "Can Large Language Models Write Parallel Code?", HPDC 2024





Introducing ParEval-Repo: Translation Tasks



Across codes chosen, using three programming model pairs:

CUDA to OpenMP Offload

Non-portable to directive-based portable

CUDA to Kokkos

Non-portable to C++ abstraction-based portable

OpenMP Threads to OpenMP Offload

CPU portable to GPU portable











Introducing ParEval-Repo: Small Cases



Starting small: three minimized codes to test capability to handle simple repos

All three include a Makefile

- I. nanoXOR: I source file
- 2. microXORh: I source file, I header file
- 3. microXOR: 2 source files, I header file

XOR kernel:





Introducing ParEval-Repo: Larger Cases



Three larger cases are taken from public codes – want to avoid codes that already have ports to minimize data contamination

- 4. **SimpleMOC-kernel**: SimpleMOC proxy, neutron flux attenuation. Has external dependency on cuRAND.
- 5. **XSBench**: OpenMC proxy, macroscopic cross-section lookup. Has public ports to OpenMP and Kokkos will show impact of possible data contamination
- 6. **Ilm.c**: CUDA implementation of LLM pretraining. We reduced the repo complexity slightly.





ParEval-Repo Translation Benchmark Suite



App. name	# of source lines	Cyclomatic complexity	# of files	OpenMP Offload	Kokkos
nanoXOR	109	33	2	×	×
microXORh	127	33	3	×	×
microXOR	133	33	4	×	×
SimpleMOC-kernel	780	59	6	×	×
XSBench	2449	264	9	✓	✓
llm.c	3039	360	7	×	×

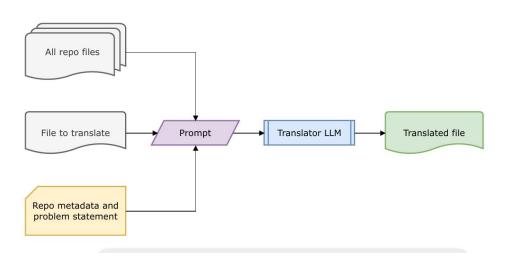




Translation Techniques: Naive



Naive translation technique: translate file-by-file in arbitrary order



Problem: exceeds LLM input context window for all but the smallest apps

PSSG

Sample prompt:

You are a helpful coding assistant. You are helping a software developer translate a codebase from the CUDA execution model to the OpenMP Offload execution model...

Below is a codebase written in the CUDA execution model... Here is the file tree of the entire repository:

<file tree>

Here is the code for each file in the codebase:

Makefile

. . .

src/main.cpp

. .

Translate the src/main.cpp file to the OpenMP offload execution model...





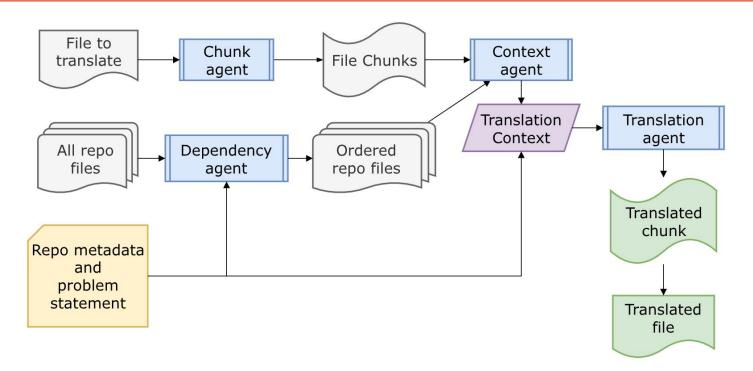
Need to break up the task into pieces that can fit into LLM context window

- Break each file into separate "chunks" that fit in context
- Order translation of files by dependency structure
 - Headers before source, source before build files)
- Use a context agent LLM to retrieve context from already-translated files for next translation





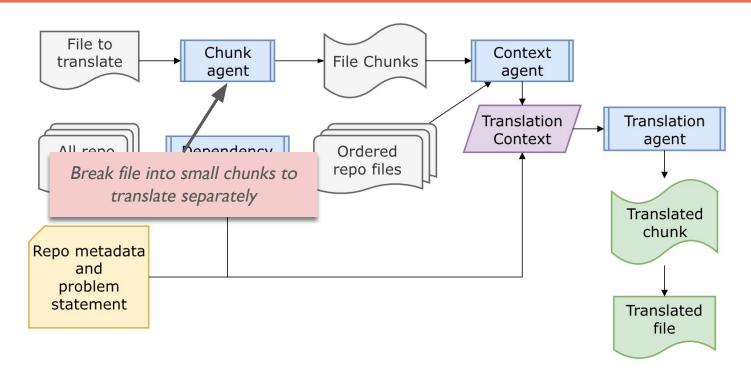








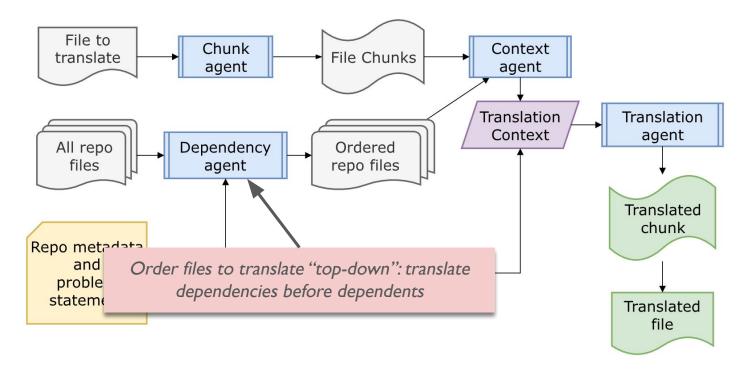








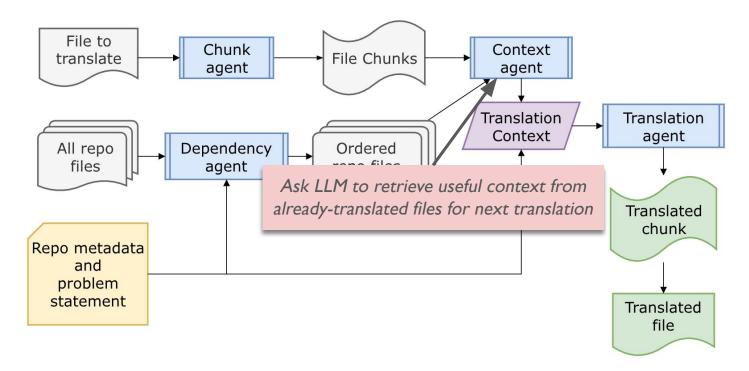








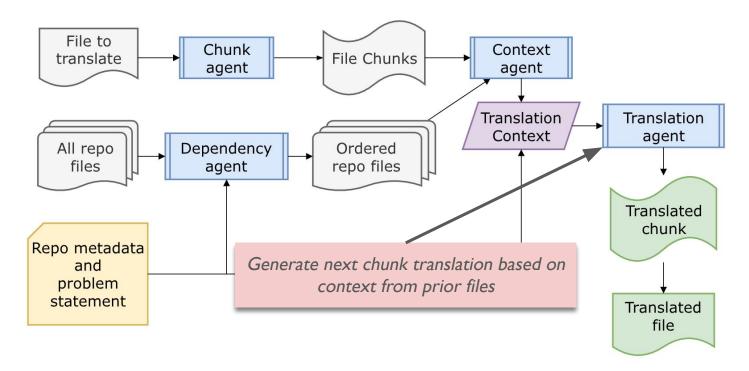








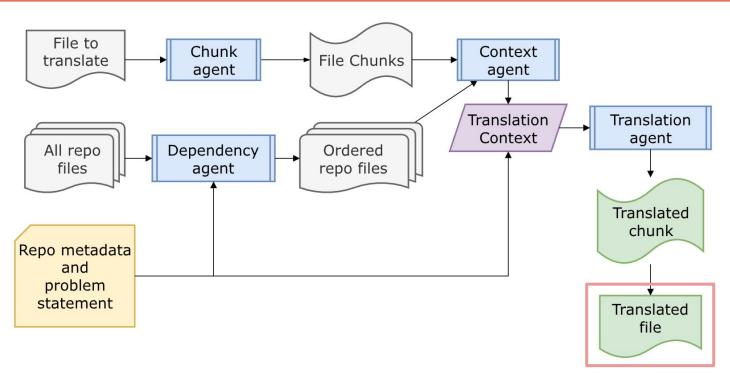
















Translation Techniques: SWE-agent





- SWE-agent is a state-of-the-art software engineering (SWE) LLM framework
- Can call tools, write tests, run code, etc.
- Needed to format our translation tasks as Github issues
- Designed for Python: editing tool cannot handle Makefiles due to tabs





Experimental Setup: LLMs Used



LLM Name	Provider	# of parameters	Open-source?	Reasoning?
Gemini 1.5 Flash	Google	?	Closed	Generic
GPT 4o-mini	OpenAl	?	Closed	Generic
o4-mini	OpenAl	?	Closed	Reasoning
Llama 3.3 70B Instruct	Meta	70 billion	Open	Generic
QwQ-32B	Alibaba Cloud	32 billion	Open	Reasoning





Experimental Setup: Evaluation Metrics



Evaluation Metrics:

- pass@ I scores: chance of generating correct translation with one attempt
 - Assessed using app's correctness test cases
- **build@** I scores: chance of generating compilable translation with one attempt

pass@ $k = \frac{1}{|T|} \sum_{p \in T} \left[1 - {N - c_t \choose k} / {N \choose k} \right]$ k = 1Number of samples generated per task k = 1Number of correct samples for task t





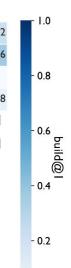
Build@I for CUDA to OpenMP Offload Translation





	Naive				
nanoXOR	1	0.98	0.92	0.92	0.9
microXORh	0	1	0.56	0.88	0.4
microXOR	0.1	0.3	0.52	0.76	0.46
SimpleMOC-kernel	0	0	0	0	0
XSBench		0	0	0	0
llm.c			0	0	0

Top-down agentic						
1	0.98	0.96	0.68	0.22		
0.24	0.24	0.12	0.36	0.36		
0	0.08	0.2	0.3	0		
0	0	0	0.02	0.08		
0	0	0	0			
0.04	0.16	0	0			



- 0.0



- Minimal success with larger codes
- When the LLM has to translate the build system as well, scores are even lower

RETRIET SHEET CHILD TOO BO CHILL SHEET CHILD TOO BO





Build@I for CUDA to OpenMP Offload Translation







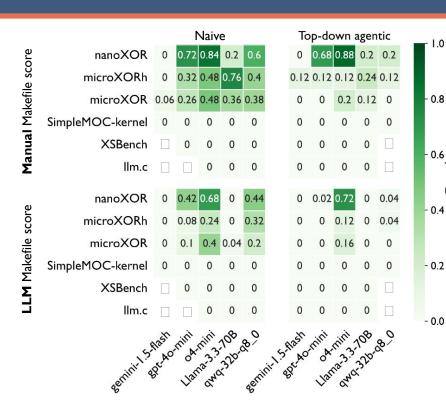
- Minimal success with larger codes
- When the LLM has to translate the build system as well, scores are even lower





Pass@I for CUDA to OpenMP Offload Translation







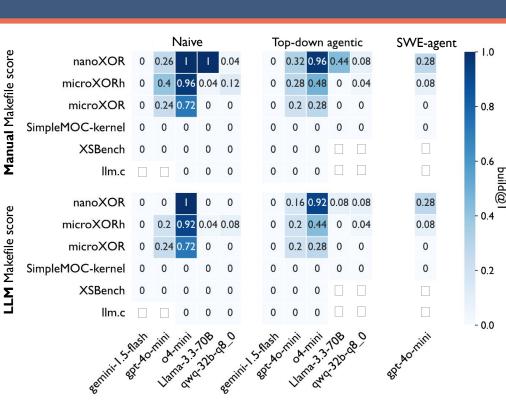
- Significant correctness cliff beyond 1-3 files
- **o4-mini naive** achieves overall best pass@l scores
- Naive often beats agentic





Build@I Results for CUDA to Kokkos







- Kokkos is even harder
 - Requires generating CMakeLists.txt
- Even with CMake provided, LLMs struggle to write compilable Kokkos
- **SWE-agent underperforms**





Error Clustering Across Translation Tasks



What are LLMs struggling with in generating buildable translations?

- Conducted clustering analysis of build outputs
- Embedded outputs to vectors with word2vec
- Clustered vectors using DBSCAN with manually tuned hyperparameters
- Manually reviewed clusters to set names, merge and split some clusters as needed

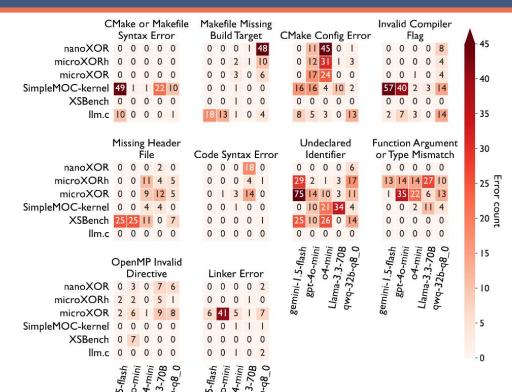




Error Clustering Across Translation Tasks



What are LLMs struggling with in generating buildable translations?



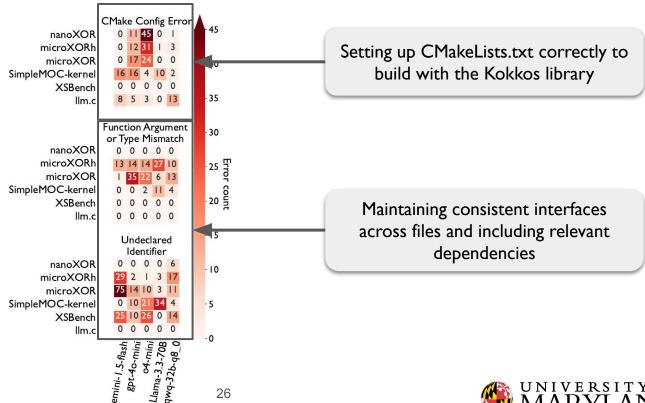




Error Clustering Across Translation Tasks



What are LLMs struggling with in generating buildable translations?







Translation Cost Analysis



Method and	nanoXOR	microXORh	microXOR
Naive o4-mini	\$0.03	\$0.03	\$0.05
Naive Llama-3.3	36 node-mins.	4 node-mins.	5 node-mins.

• LLM translation is expensive

We derive expected tokens used with

Average token cost per generation for the task $E_{\kappa} = \left(\frac{1}{\text{pass@1}}\right) \cdot \frac{1}{\kappa}$ Single generation correctness chance

Convert this to node hours or API
cost (\$) as relevant for the model





Conclusion and Future Work



- Employing LLMs to automate application translation to new programming models has tremendous potential for developer productivity
- But, state-of-the-art LLMs struggle to correctly translate full applications
- Translating build systems and cross-file dependencies are most significant hurdles
- Opportunities for improvement: building a dataset for fine-tuning or improving prompt context, or building a more sophisticated agent approach



← ParEval-Repo is on GitHub!

Contact: jhdavis@umd.edu





jhdavis@umd.edu





ParEval-Repo Translation Benchmark Suite

App. name	# of source code lines	Cyclomatic complexity	# of files	OpenMP Threads	OpenMP Offload	CUDA	Kokkos
nanoXOR	109	33	2	✓	?	✓	?
microXORh	127	33	3	✓	?	✓	?
microXOR	133	33	4	✓	?	1	?
SimpleMOC-kernel	780	59	6		?	✓	?
XSBench	2449	264	9	1	√?	✓	√?
llm.c	3039	360	7		?	✓	?
		Port alı	ready exists	Port doe	sn't exist,		



will translate

