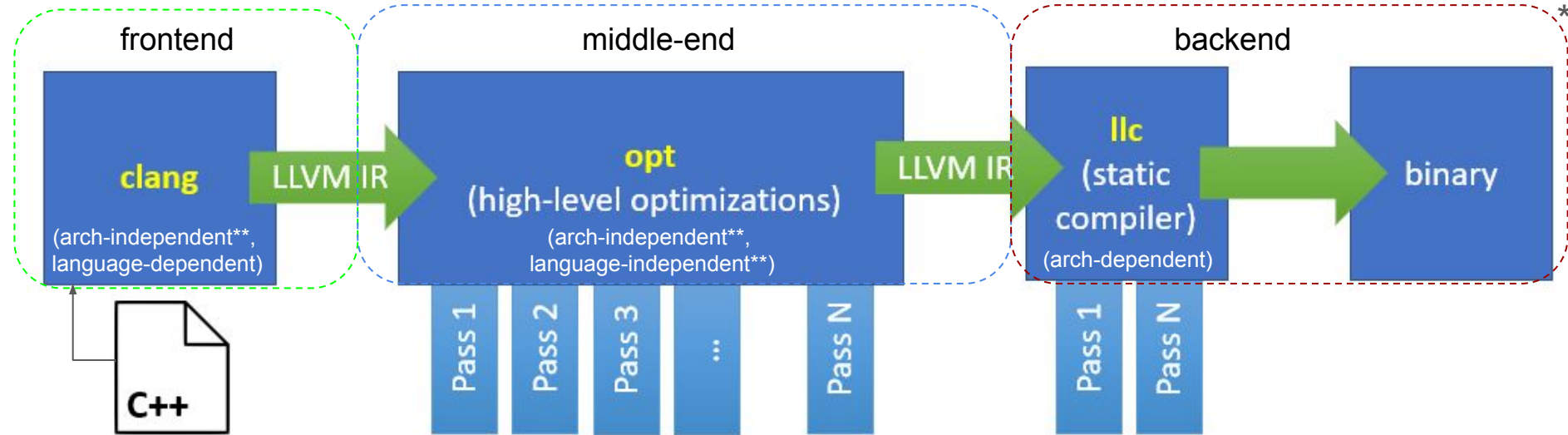


Compiler Auto-tuning: Automatic search of optimal pass sequences / parameters by ML



Nikolay Efanov, PhD, associate professor at MIPT
Aleksy Otrashchenko, student at MIPT

* Clang / LLVM
** There are some nuances

Compiler auto-tuning:

Compilers implements a lot of different transformation passes

- LLVM: ~160
- GCC: ~250
- At all, there about ~500 different types of passes
- Some passes are parametric (“loop-unroll”, etc.)

Pass sequences, which are “optimal in mean” for different optimization targets (size, perf, calculation accuracy), are grouped into “optimization levels”:

- O0, O1
- O2 (usually used in industry)
- O3
- Os, Oz...

Passes have different granularity levels, and require different program representations

Some programs (mostly all, in fact) can be better optimized, than -O*, if to tune the passes and their parameters more accurately.

The sequence of passes also oftenly important (AB != BA).

Problems and combinatorial estimations of search spaces

- Choosing the right set of phases:

$$|\Omega_{selection}| = \{0, 1\}^n \quad (1)$$

(1) Estimates the search space size for binary selection of phases. For example, if $n = 9$ then $|\Omega_{selection}| = 2^9 = 512$. Moreover, if taking into account parameterization of phases and phases repetition (some of phases can be applied more than 1 time), the extended estimation:

$$|\Omega_{selection_extended}| = \prod_{j=1}^{n'} \{0, 1, \dots, m_j\}, \text{ where} \quad (2)$$

$m_j + 1$ is the total number of parameters' values for j -th phase, $j \in [1, n']$, and n' is number of phases to apply.

- Phase-ordering problem: Due to the permutations, the search space size of this problem growth as factorial (3):

$$|\Omega_{phases}| = n! \quad (3)$$

Moreover, taking into account possible repetitions of phases and variability of sequence length, this estimation becomes (4):

$$|\Omega_{phases_extended}| = \sum_{i=0}^l n^i, \text{ where} \quad (4)$$

l is the max length of phases sequence vector.

Thus, there is, for example, 3905 variants for $l=5$

Note 1:

Methods for automation of phases choosing, phase-ordering and phase parameters picking required

Note 2:

There is no "ideal" order of passes applying in common case. Pass A transforms code in that it can corrupt the optimizations, which could be successfully performed by the pass B next to the A.

Compiler auto-tuning: passes reordering

Different objectives, pass sets, even constrains:

Objective function
Constrains

For whole program p granularity

$$size(C(s, p)) \Delta_{rt}(C(s, p), T, O_2) < \epsilon \rightarrow \min, \text{ где}$$

For function-level granularity

$$\sum_{i=1}^{|F(\epsilon)|} \widehat{size}(C(s_i, f_i)) \Delta_{rt}(C(s_i, f_i), T, O_2) < \epsilon \rightarrow \min$$

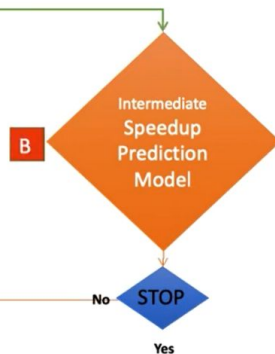
Different solution approaches:



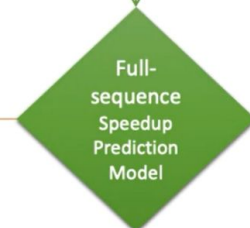
MDP* (POMDP)



Prog. Features Opt. Passes



Classification



1. Iterative (n passes are applied, predict n+1)

- + Models can be relative simple
- One characterization per step
- Local min convergency risks

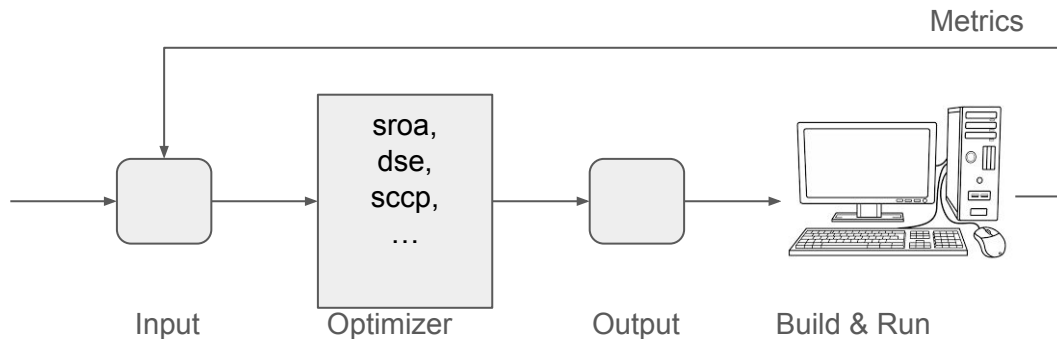
2. Full-sequence prediction

- Models are more complicated
- Requires accurate program representations
- + Requires only one characterization on inference
- + Much resistance to local min

- **Support by compiler**
 - Switching passes on / off (in fixed order)
 - Pass-reordering
 - Infrastructure for working with IR, representations, optimizer, etc
- **Optimal sequences search**
 - Optimization space exploration
 - Benchmarks preparation & analysis
 - Algorithm choosing
 - EA, RL
 - Supervised learning, collaborative filtering
- **Code characterization**
 - By explicit (expert) features-selection
 - Graph-structured representation analysis (flow-aware)
 - By ML
- **Specialization for target platform**
 - Parameterized passes, parameters choosing
 - Prior knowledge about hardware features

Iterative compilation [1]

Sequential code re-compilation with different pass sequences for choosing the best (k-best)



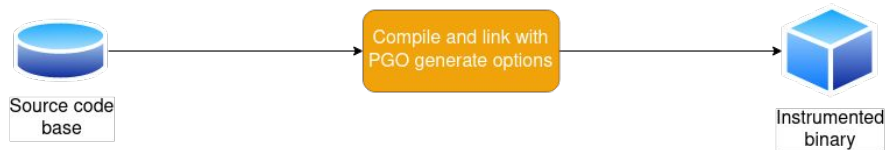
- + Guaranteed best results via brute-force
- + Very simple approach
- Monstrous overhead / time consumption
- No knowledge transfer between experiments

1. Bodin, François & Kisuki, Toru & Knijnenburg, Peter & Boyle, Mike & Rohou, Erven. (2000). Iterative compilation in a non-linear optimisation space. Workshop on Profile and Feedback-Directed Compilation.

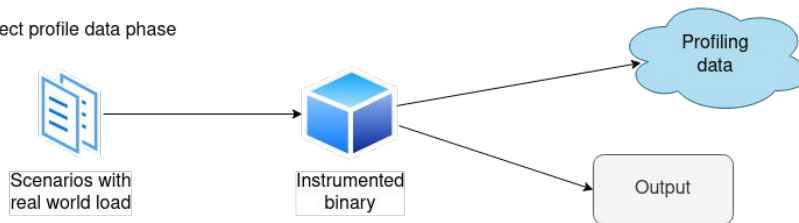
Adaptive compilation (“Profiling guided optimization”) [2]

Profiling of instrumented program and applying of the best passes for collected profile

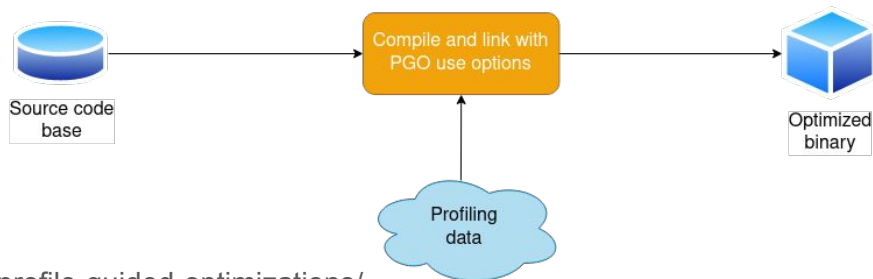
Instrument phase



Collect profile data phase



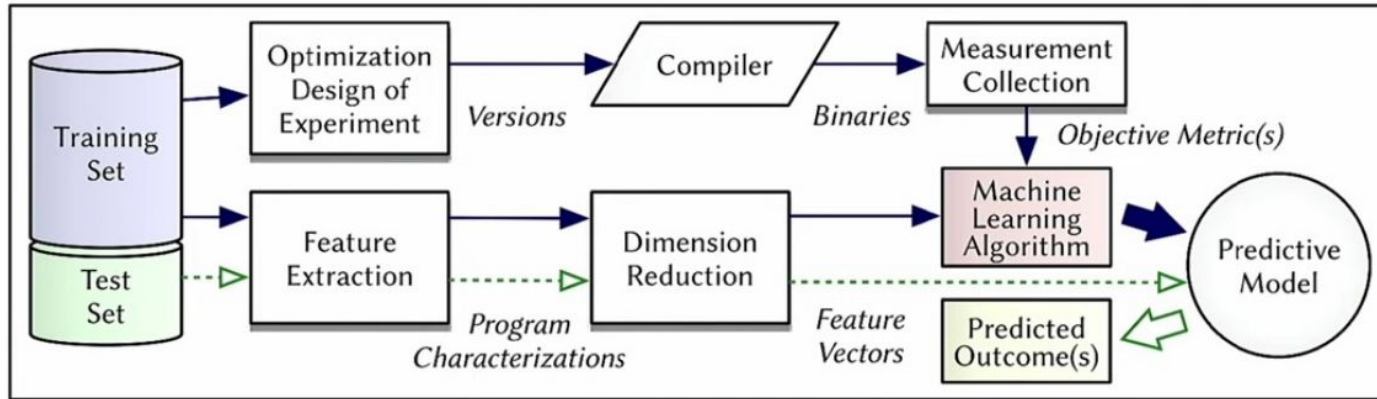
Optimize phase



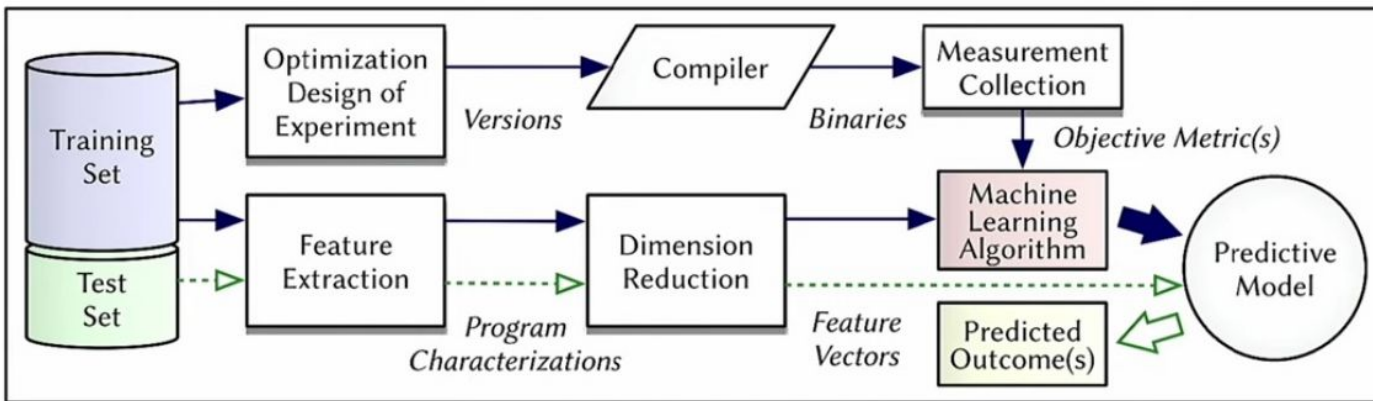
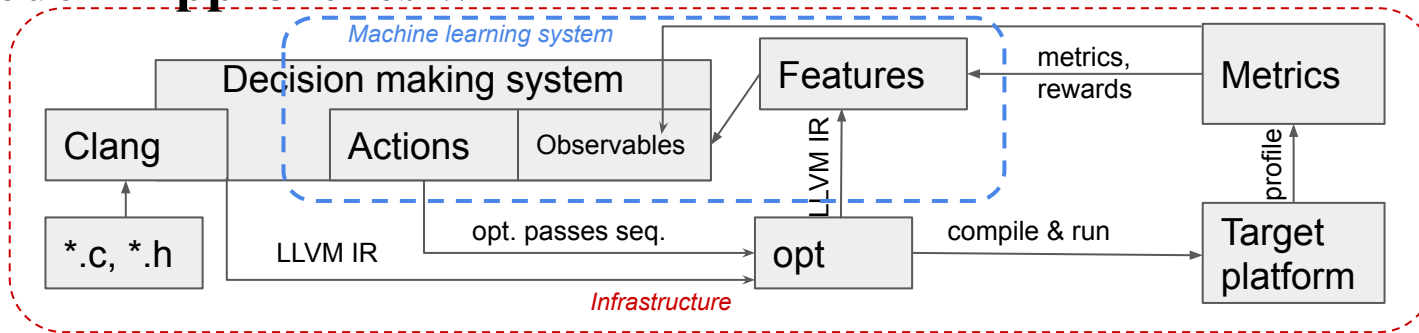
- + Relatively simple approach
- + Approbated in industry (Firefox Mozilla, etc)
- No knowledge transfer
- Best results for certain profiles (and may be worse for another scenarios)

Modern approaches with ML

- + Knowledge transfer
- + Generalization
- Quite complex
- Require novel code characterization methods, benchmarks preparation



Modern approaches with ML



Ускорение поиска и трансфера знаний

Итеративная компиляция заложила основу автоматизированного подхода [Bodin и др., 1998-2000] и определила дальнейший вектор его развития. За ~25 лет успехи следующие:

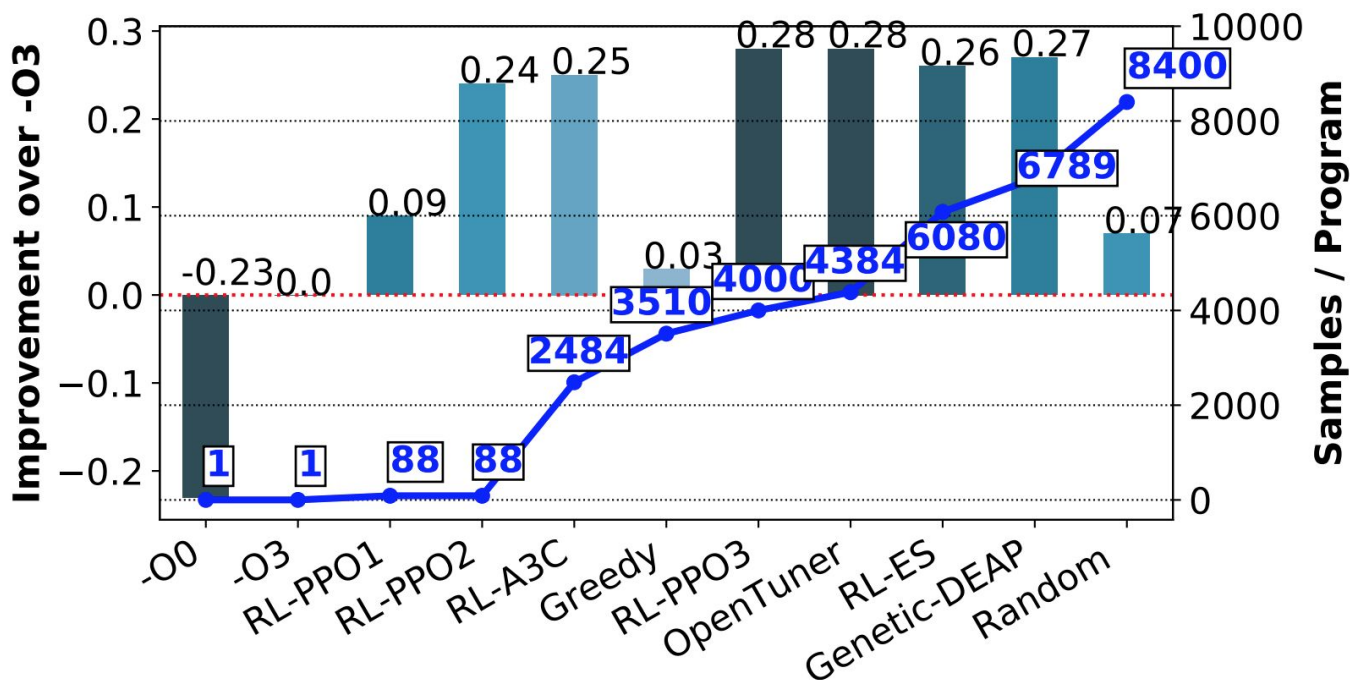
- Последние десятилетия ускорение поиска осуществлялось в основном за счёт подходов с ГА и стохастический поиск [Tagtekin и др., 2021], а также комбинированием ГА с обучением с учителем [Agakov et al, 2006], etc.
- В последние 5-10 лет начали применять обучение с подкреплением, практические результаты демонстрируют существенно более быструю сходимость к результатам, как у ГА. При этом возникает возможность трансфера знаний
- Задан тренд на построение сложных методов характеристики, в т.ч. посредством обучения представлений
- Разработаны решения, использующие информацию о hardware для лучшей оптимизации (NeuroVectorizer)
- Исследовательские работы, как правило, не связываются со сложными критериями оптимизации (однокритериальная с ограничениями, многокритериальная). В основном, это победа над -O3 / -Os /-Oz
- Наша группа, наоборот, использует нецелевые характеристики программ как ограничения
- Решения в данной технической области крайне быстро устаревают

Bodin, François & Kisuki, Toru & Knijnenburg, Peter & Boyle, Mike & Rohou, Erven. (2000). Iterative compilation in a non-linear optimisation space. Workshop on Profile and Feedback-Directed Compilation.

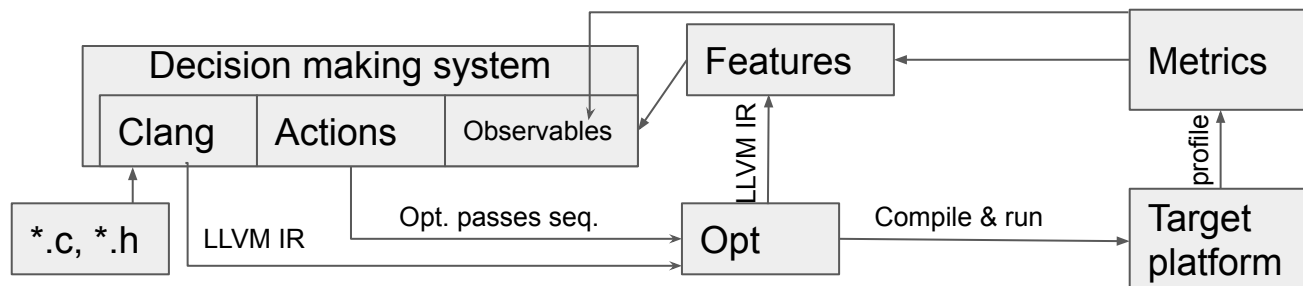
Tagtekin, B., Höke, B., Sezer, M. K., & Öztürk, M. U. (2021, August). FOGA: Flag Optimization with Genetic Algorithm. In 2021 International Conference on INnovations in Intelligent Systems and Applications (INISTA) (pp. 1-6). IEEE.

F. Agakov, E. Bonilla, J. Cavazos, B. Franke, G. Fursin, M. F. O'Boyle, J. Thomson, M. Toussaint, and C. K. Williams, "Using machine learning to focus iterative optimization," in International Symposium on Code Generation and Optimization (CGO'06). IEEE, 2006, pp. 11–pp.

Appendix #1



Appendix #2: Proof-of-concept with heuristic search, LLVM [Efanov, 2023]



Results for CBench .TEXT size reduction without runtime degradation:

benchmark	baseline	Best result	Gain, %
patricia			23.356009070294785
gsm			21.191308866697568
blowfish			10.062516626762436
bzip2			7.324845763802766
jpeg-d			7.229062428095125
Jpeg-c			7.150153217568948
sha			6.361275964391691
tiffmedian			5.311668039060172
tiffdither			5.22980749841004
tiff2bw			5.21401899703284
tiff2rgba			5.20897960871903
qsort			4.854968113556881
stringsearch			3.0260047281323876
stringsearch2			2.7645788336933044
crc32			1.4736842105263157
dijkstra			0.7469654528478058

Method: heuristic search (least from positives per-step);

Iterations num.: 100 Iterations;

Episode len.: 15;

Patience: 5;

Runtime eval: 10 times, mean;

Env.: LLVM;

Benchmark: CBench-v1;

Actions: Oz_extra;

Target: size reduction under const performance degradation

Challenges:

- Off-line full sequence prediction (currently solved, can be improved)
- Static code characterization (flow-aware & scalable) methods construction & integration
- Convergence speed-up by subsequences extracted from ODG
- Optimal parameters prediction for parameterized passes (loop unroll, vectorize, inline, etc)
- ML methods improvement (achieved ~11% max size reduction on CBench with AC RL)

Place for slides from En&T-2023 Conference...

Code characterization

Code characterization

- **Static**
 - Aggregated features (for ex, extracted by NLP [3] models or collected from IR by values of instructions, BBs, functions, phi-nodes counters, etc [1,4-5].)
 - Graph-structured (CFG, DFG, CallGraph, etc [4])
 - Mixed
- **Dynamic**
 - Arch-dependent (perf counters, etc) [2-4]
 - Arch-independent (for ex, number of function calls)
- **Hybrid**
 - Collect as much as ...

1. Q. Huang, et al., "AutoPhase: Compiler Phase-Ordering for HLS with Deep Reinforcement Learning," in 2019 IEEE 27th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM), San Diego, CA, USA, 2019 pp. 308-308.

2. POSET-RL: <https://arxiv.org/abs/2208.04238>

3. [Wang, H, Tang, Z, Zhang, C et al]. (4 more authors) (2022) Automating Reinforcement Learning Architecture Design for Code Optimization. In: CC 2022: Proceedings of the 31st ACM SIGPLAN International Conference on Compiler Construction. The 31st ACM SIGPLAN International Conference on Compiler Construction (CC '22), 02-03 Apr 2022.

4. S. VenkataKeerthy, Rohit Aggarwal, Shalini Jain, Maunendra Sankar Desarkar, Ramakrishna Upadrasta, and Y. N. Srikant. 2020. IR2VEC: LLVM IR Based Scalable Program Embeddings. ACM Trans. Archit. Code Optim. 17, 4, Article 32 (December 2020), 27 pages.

5. Chris Cummins, Bram Wasti, Jiadong Guo, Brandon Cui, Jason Ansel, Sahir Gomez, Somya Jain, Jia Liu, Olivier Teytaud, Benoit Steiner, Yuandong Tian, and Hugh Leather. 2022. CompilerGym: robust, performant compiler optimization environments for AI research. In Proceedings of the 20th IEEE/ACM International Symposium on Code Generation and Optimization (CGO '22). IEEE Press, 92–105.

Code characterization

- **Static**
 - Aggregated features (for ex, extracted by NLP [3] models or collected from IR by values of instructions, BBs, functions, phi-nodes counters, etc [1,4-5].)
 - Graph-structured (CFG, DFG, CallGraph, etc [4])
 - Mixed
- **Dynamic**
 - Arch-dependent (perf counters, etc) [2-4]
 - Arch-independent (for ex, number of function calls)

Note:

- Несмотря на “раскрученность”, модели на основе NLP (code2vec, CodeBERT, etc) применяются не так часто, что обусловлено многими факторами: нечувствительностью к потоковой информации, трудностью определения семантики (и, как результат, невозможностью обеспечения некоторых свойств представления). Модели, для которых операционная семантика определена, оперируют более сложными структурами контекста, чем скип-грамма из последовательности слов. (“Код -- это не текст”)
- Решения на основе LLM появляются (первая работа -- сентябрь 2023 г.). Качество результатов -- пока что не высокое (3% для уменьшения количества инструкций без ограничений на LLVM).

<https://arxiv.org/pdf/2309.07062.pdf>

Code characterization: from experience view

- NLP-based models are inaccurate (program is not a token sequence, relations between entities should be extracted) [1]. Relations and semantics matters [4,11].
 - Thus, graph-based characterization methods required (or another kind, which take into account the relations on instructions, variables, arguments, types, etc)
 - Theoretically, inst2vec should give the best results in comparison with another, Because of taking into account graph-structure of program flows. But it is not enough.
- Manually-crafted features (instr.number, bb sizes, trip-counts) are specialization (suitable for certain applications rather than for general purpose)
 - Question of automation is opened -- defaultly, need experts
- Dimensionality reduction required
 - In case IR2Vec, the inference sometimes taken in minutes for relative big programs (for ~600 lines in TU)
 - Sparsity reduction of feature space will lead, at least, to convergence speed gain
- Models
 - Should keep of semantics (demonstrated [11], that CFG,DFG,Call,Type graph-based methods are relatively good in this criteria).
 - Should be extensible to represent static & dynamic features
- The first model with LLMs [<https://arxiv.org/pdf/2309.07062.pdf>] shows some minor results for unconstrained size reduction.

Code characterization: AutoPhase repr. Vs InstCount

AutoPhase

TotalInsts	567
TotalMemInst	340
testUnary	244
const32Bit	133
NumLoadInst	127
NumEdges	95
NumStoreInst	91
TotalBlocks	73
BBNoPhi	73
NumCallInst	66
NumBrInst	66
BranchCount	66
BlockLow	66
const64Bit	65
numConstOnes	59
NumBitCastInst	56
numConstZeroes	52
onePred	45
oneSuccessor	40

...

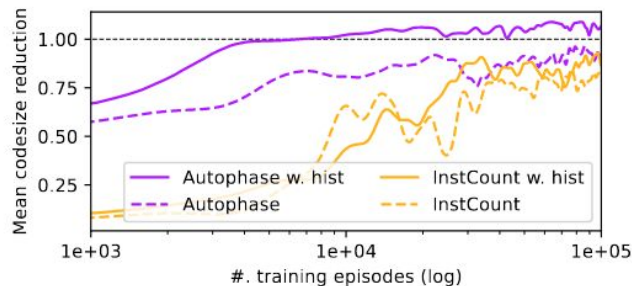
InstCount

TotalInstsCount	567
LoadCount	127
StoreCount	91
TotalBlocksCount	73
CallCount	66
BrCount	66
BitCastCount	56
AllocaCount	37
ICmpCount	29
GetElementPtrCount	19
TotalFuncsCount	14
SExtCount	13
AddCount	13
AndCount	9
TruncCount	8
AShrCount	7
SubCount	6

Accurate choosing of representation

leads to:

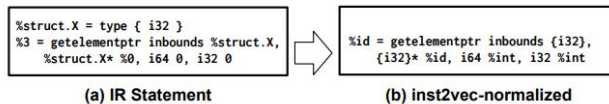
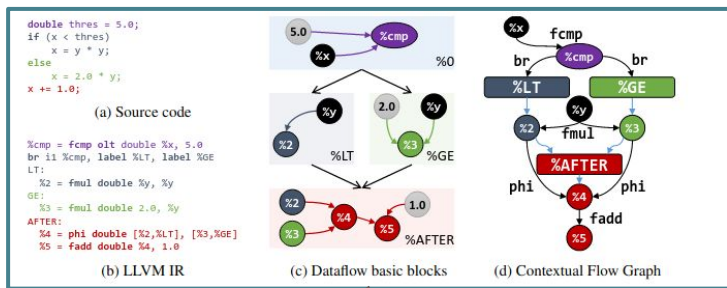
- Convergence speedup
- Results improvement



In both cases stronger performance is achieved when coupling the program representation with a histogram of the agent's previous actions. The AutoPhase representation encodes more attributes of the structure of programs than InstCount and achieves greater performance

Характеризация кода: Inst2vec (2018)

“Statements that occur in the same contexts tend to have similar semantics.”



Similarity To define similarity, one first needs to define the *semantics* of a statement. We draw the definition of semantics from Operational Semantics in programming language theory, which refers to the effects (e.g., preconditions, postconditions) of each computational step in a given program. In this paper, we specifically assume that each statement modifies the system state in a certain way (e.g., adds two numbers) and consumes resources (e.g., uses registers and floating-point units). It follows that semantic similarity can be defined by two statements consuming the same resources or modifying the system state in a similar way. Using this definition, two versions of the same algorithm with different variable types would be synonymous.

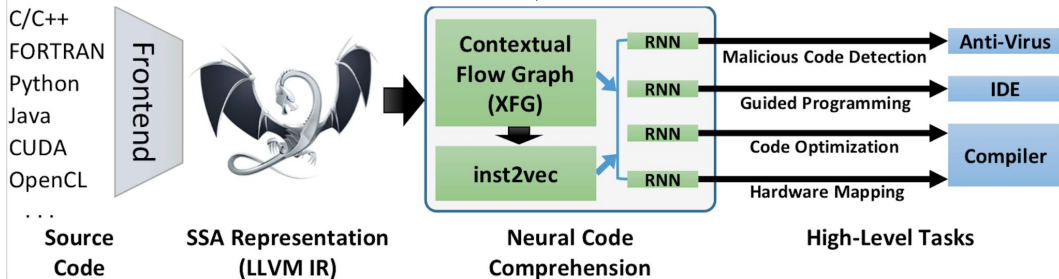


Table 2: Analogy and test scores for *inst2vec*

Context type	Context Size	Syntactic Analogies		Semantic Analogies		Semantic Distance Test
		Types	Options	Conversions	Data Structures	
CFG	1	0 (0%)	1 (1.89%)	1 (0.07%)	0 (0%)	51.59%
	2	1 (0.18%)	1 (1.89%)	0 (0%)	0 (0%)	50.47%
	3	0 (0%)	1 (1.89%)	4 (0.27%)	0 (0%)	53.79%
DFG	1	53 (9.46%)	12 (22.64%)	2 (0.13%)	4 (50.00%)	56.79%
	2	71 (12.68%)	12 (22.64%)	12 (0.80%)	3 (37.50%)	57.44%
	3	67 (22.32%)	18 (33.96%)	40 (2.65%)	4 (50.00%)	60.38%
XFG	1	101 (18.04%)	13 (24.53%)	100 (6.63%)	3 (37.50%)	60.98%
	2	226 (40.36%)	45 (84.91%)	134 (8.89%)	7 (87.50%)	79.12%
	3	125 (22.32%)	24 (45.28%)	48 (3.18%)	7 (87.50%)	62.56%

Table 3: Algorithm classification test accuracy

Metric	Surface Features [49] (RBF SVM + Bag-of-Trees)	RNN [49]	TBCNN [49]	inst2vec
Test Accuracy [%]	88.2	84.8	94.0	94.83

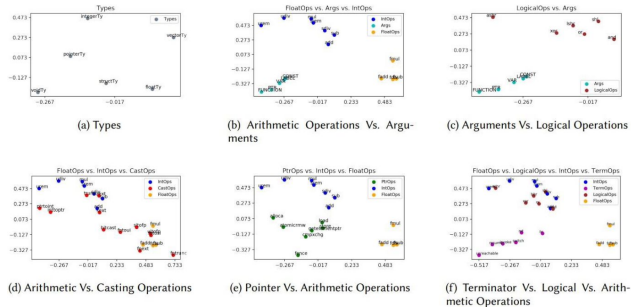
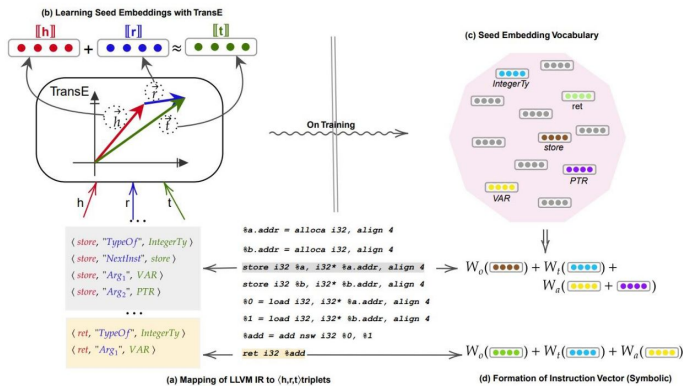
1. Read LLVM IR statements once, storing function names and return statements.
2. Second pass over the statements, adding nodes and edges according to the following rule-set:
 - (a) Data dependencies within a basic block are connected. Note: that's way data-flow is taken into account
 - (b) Inter-block dependencies (e.g., ϕ -expressions) are both connected directly and through the label identifier (statement-less edges).
 - (c) Identifiers without a dataflow parent are connected to their root (label or program root).
 - (d) Calls to external code (e.g., libraries, frameworks) are divided into two categories: statically-(connections) and dynamically-linked (stubs).
3. Achieved XFG becomes a context for *inst2vec* statement (упрощённый LLVM IR). Then, skip-gram learning is used to learn the embedding (XFG as set of paths)

Code characterization: IR2Vec (2020)

Seed dictionary (pre-trained) for LLVM IR entities + combination rules

Lets consider triplets $\langle h, r, t \rangle$, where entities h, t have relation r of one of three types:

1. *TypeOf* – relation between instruction code and current instruction type
2. *NextInst* – relation between current and next instruction code
3. *Arg_i* – between instruction code and its i -th operand



Lets consider the entities, which define the instruction I : code $O^{(I)}$, i -th argument $A_i^{(I)}$, $i \in [0, n]$, type $T^{(I)}$ with vector representations $[O^{(I)}], [A_i^{(I)}], [T^{(I)}]$. Then, the instruction $\langle O^{(I)}, T^{(I)}, A_0^{(I)}, \dots, A_n^{(I)} \rangle$ can be represented as vector :

$$W_o[O^{(I)}] + W_t[T^{(I)}] + W_a\left(\sum_{i=0}^n [A_i^{(I)}]\right), \text{ where} \quad (1)$$

$W_o > W_t > W_a$ – weight coefficients, learned by TransE method.

Let RD_0, \dots, RD_m are reaching-definitions of some argument $A_j^{(I)}$ and its vector representations are $[RD_0], \dots, [RD_m]$. Then, $A_j^{(I)}$ can be represented as:

$$[A_j^{(I)}] = \sum [RD_i^{(I)}] \quad (2)$$

Note: that's way data-flow is taken into account

Next, for each basic block BB_j representation can be calculated as sum of live instructions LI_0, \dots, LI_k representations:

$$[BB_j] = \sum_{i=0}^k [LI_i] \quad (3)$$

Finally, function F can be represented as sum of its basic blocks BB_0, \dots, BB_b representations:

$$[F_j] = \sum [BB_i] \quad (4)$$

Note: that's way control-flow is taken into account

In the same way, the vector of whole translation unit P can be calculated as sum of vector representation of functions F_1, \dots, F_f :

$$[P] = \sum_{i=0}^f [F_i] \quad (5)$$

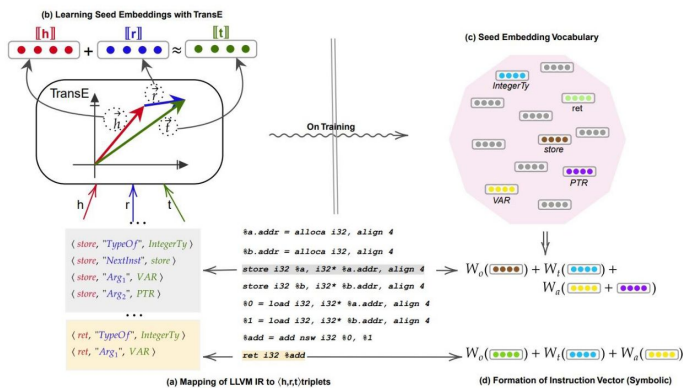
Thus, hierarchical representation method, which can map instructions, types, arguments and relations between them, basic blocks, functions and whole program is constructed.

Code characterization: IR2Vec (2020)

Seed dictionary (pre-trained) for LLVM IR entities + combination rules

Lets consider triplets $\langle h, r, t \rangle$, where entities h, t have relation r of one of three types:

1. *TypeOf* – relation between instruction code and current instruction type
2. *NextInst* – relation between current and next instruction code
3. *Arg_i* – between instruction code and its i -th operand



Lets consider the entities, which define the instruction I : code $O^{(I)}$, i -th argument $A_i^{(I)}$, $i \in [0, n]$, type $T^{(I)}$ with vector representations $[O^{(I)}], [A_i^{(I)}], [T^{(I)}]$. Then, the instruction $\langle O^{(I)}, T^{(I)}, A_0^{(I)}, \dots, A_n^{(I)} \rangle$ can be represented as vector :

$$W_o[O^{(I)}] + W_t[T^{(I)}] + W_a\left(\sum_{i=0}^n [A_i^{(I)}]\right), \text{ where} \quad (1)$$

$W_o > W_t > W_a$ – weight coefficients, learned by TransE method.

Let RD_0, \dots, RD_m are reaching-definitions of some argument $A_j^{(I)}$ and its vector representations are $[RD_0], \dots, [RD_m]$. Then, $A_j^{(I)}$ can be represented as:

$$[A_j^{(I)}] = \sum_{i=0}^m [RD_i^{(I)}] \quad (2)$$

Note: that's way data-flow is taken into account

Next, for each basic block BB_j representation can be calculated as sum of live instructions LI_0, \dots, LI_k representations:

$$[BB_j] = \sum_{i=0}^k [LI_i] \quad (3)$$

Finally, function F can be represented as sum of its basic blocks BB_0, \dots, BB_b representations:

Note:

- Data-flow sensitive
- Questions to CFG sensitivity (“+” commutativity)
- Questions to original implementation (NextInst is not used)...
- Previously approved by LLVM-community for MLGO integration
- Quite strange but practically applicable

Code characterization: IR2Vec [4,8] : Experiments & Alternatives

TransE pre-trained LLVM IR entities+manual rules

In [Zavodskikh et al, 2022] Zavodskikh R.K. (my PhD student) have experimentally checked the ability of IR2Vec to save statically estimated reuse-distance between accessed array elements for chosen types of loops. It is determined that the method reflects the instrumentation data in representation vectors, and results of estimation are correlated with Linux Perf measurements. Thus, it demonstrates the ability of IR2Vec to represent control-flow and data-flow accurately enough.

Nevertheless, there are a list of drawbacks of the method detected:

- BBs are not ordered according to the original CFG order.
- Original implementation of the method not takes into account NextInst relations (unlike in the method description [4]).

Thus, the implementation should be improved.

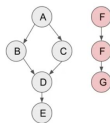
Also, there are alternative methods of flow-aware program features extraction recently published (<https://ieeexplore.ieee.org/document/9275317>, <https://chriscummins.cc/2017/deep-learning-in-compilers>), mostly based on graph-structures deep learning via message-passing graph neural networks. But this solutions seems more hard to implement (except PrograML [11], which already integrated to CompilerGym), re-train and use, in contrast with IR2Vec, in which the basic relations on IR entities are extracted manually and pre-trained as a seeds for representation of compound entities as superpositions of the constituents vectors.

Code characterization: ProGraML (2021)

Graph-structured, flow-aware, type-aware representation, forwarded to MP-GNN

It isn't a vector of numbers

Feature vectors are easy to fool (e.g. insert **dead code**).



It isn't a sequence of tokens

Sequential representations fail on non-linear relations, **long-range deps**.

```
void A(int a)
int b = init();
// ... 1000 lines
// ...
return b - a;
}
```

Learning with ProGraML: GGNNs

Message Passing

$$M(h_{u}^{t-1}, e_{uv}) = W_{\text{type}(e_{uv})} (h_u^{t-1} \odot p(e_{uv})) + b_{\text{type}(e_{uv})}$$

4 typed weight matrices for (forwards/backwards) (control/data/call) edge types

Readout Head

$$R_i(h_i^T, h_i^B) = \sigma(f(h_i^T, h_i^B)) \cdot g(h_i^T)$$

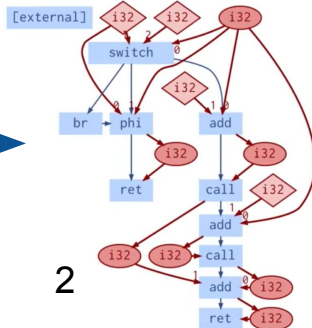
Position going to different control branches and opened/closed

per-vertex prediction after T message-passing steps

Add graph vertices for **constants** (diamonds) and **variables** (oblongs).

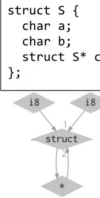
Edges are **data-flow**.

Edge position attribute for **operand order**.



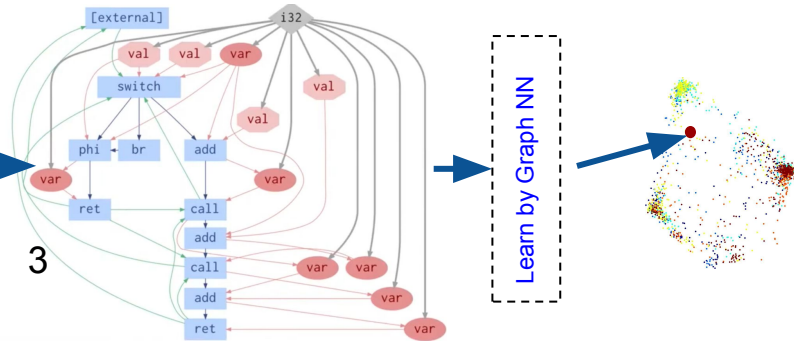
Nodes represent **types**. Edges are **instances**.

Types are **composable**. Edge position per field.



Embedding construction approach:

- inst2vec for nodes-statements
- Relations learning: MP-GNN
- In some sense, this is mixed approach



Learn by Graph NN

Note: Data-Flow

Note: CallGraph and TypeInfo*

Deep Data Flow

Dataset: 40k LLVM IRs covering 5 programming languages

	inst2vec	CDFG	ProGraML
Reachability Trivial forwards control-flow E.g. dead code elimination	0.012	0.998	0.998
Dominance Forwards control-flow E.g. global code motion	0.004	0.999	1.000
Data Dependencies Forwards data-flow E.g. instruction selection	-	-	0.997
Live-out Variables Backwards control- and data-flow E.g. register allocation	-	-	0.937
Global Common Subexpressions Instruction/operand sensitive E.g. GCS Elimination	0.000	0.009	0.996

Use vertex labels as embedding keys



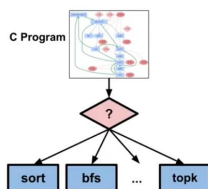
Derive vocab from set of unique vertex labels on **training graphs**.

Separate type/instruction nodes leads to **compact vocab**, excellent coverage on unseen programs compared to prior approaches:

	Vocabulary size	Test coverage
inst2vec [12]	8,565	34.0%
CDFG [14]	75	47.5%
ProGraML	2,230	98.3% *without types

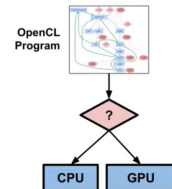
inst2vec: combined instruction+operands `i32 <id> = <id> <int8>`
 CDFG: uses only instructions for vocab, ignores data

1. Algorithm Classification



1.35x improvement over state-of-art

2. Heterogeneous Device Mapping



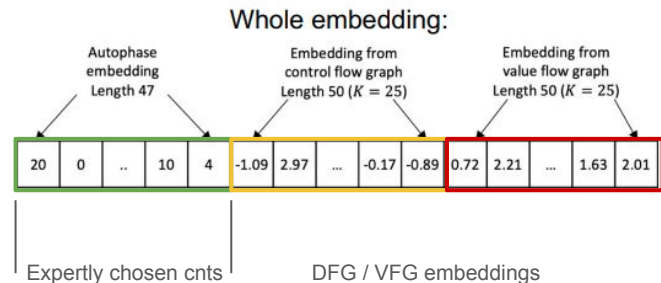
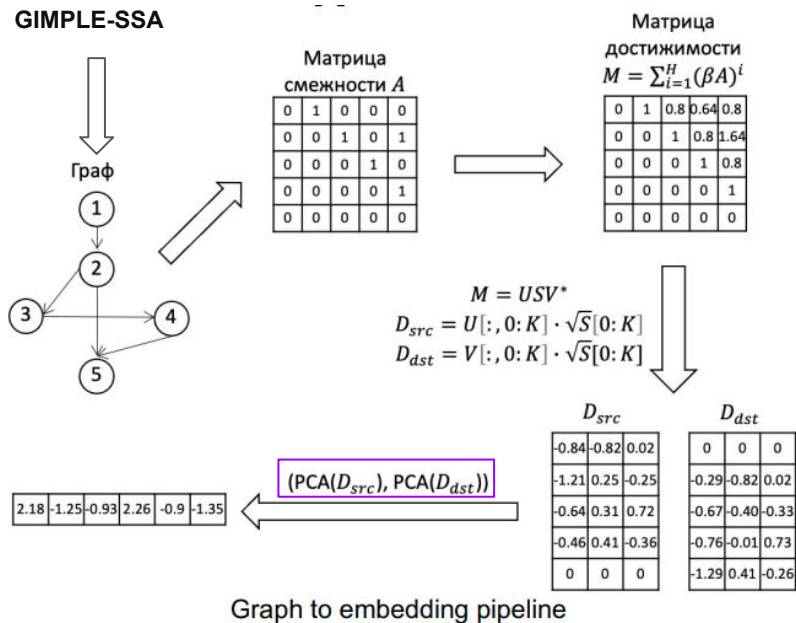
1.20x improvement over state-of-art

* There some type elision in LLVM IR

Code characterization: mixed static methods (static features + CFG / VFG)

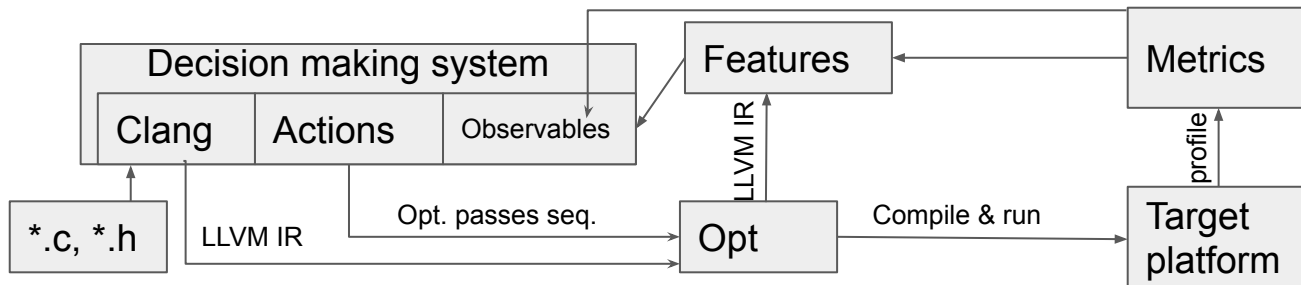
GCC GIMPLE-SSA embedding [Otrashchenko, Akimov, Efanov, 2023]

- GCC IR is characterised using autophase characterisation, control and value flow graphs
- Autophase characterisation consists of information about IR, available immediately during compilation
- The embedding from control flow graph and value flow graph are acquired as shown on the picture



RL/GA Approach: GCC PHO (external ppt)

Appendix 1: Proof-of-concept with heuristic search



Results for CBench .TEXT size reduction without runtime degradation:

benchmark	baseline	Best result	Gain, %
patricia			23.356009070294785
gsm			21.191308866697568
blowfish			10.062516626762436
bzip2			7.324845763802766
jpeg-d			7.229062428095125
Jpeg-c			7.150153217568948
sha			6.361275964391691
tiffmedian			5.311668039060172
tiffdither			5.22980749841004
tiff2bw			5.21401899703284
tiff2rgba			5.20897960871903
qsort			4.854968113556881
stringsearch			3.0260047281323876
stringsearch2			2.7645788336933044
crc32			1.4736842105263157
dijkstra			0.7469654528478058

Method: heuristic search (least from positives per-step);

Iterations num.: 100 Iterations;

Episode len.: 15;

Patience: 5;

Runtime eval: 10 times, mean;

Env.: LLVM;

Benchmark: CBench-v1;

Actions: Oz_extra.

Challenges:

- Off-line full sequence prediction
- Static code characterization (flow-aware & scalable) methods construction & integration
- Convergence speed-up by subsequences extracted from ODG
- Optimal parameters prediction for parameterized passes (loop unroll, vectorize, inline, etc)
- ML methods improvement (achieved ~11% max size reduction on CBench with AC RL)

3. Automating Reinforcement Learning Architecture Design for Code Optimization (2022)

```

1 import SuperSonic as ss
2 from SuperSonic.statefunctions.models import *
3 ...
4 statefs = [Word2Vec(...), Doc2Vec(...), CodeBERT
5 (...), ActionHistory(...)]
6 tranfs = [DNN(...), CNN(...), LSTM(...)]
7 rewards = [RelativeMeasure(...), tanh(...)]
8 rl_algs = [MCTS(...), PPO(...), DQN(...), QLearning
9 (...)]
10 actions = [Init(...)]
11 ...
12 class SuperOptimizer(ss.PolicyInt):
13     def __init__(self, statefs, tranfs, actions,
14                 rewards, rl_algs):
15         self.PolicySpace = {
16             "StatList": statefs,
17             "TranList": tranfs,
18             "ActList": actions,
19             "RewList": rewards,
20             "AlgList": rl_algs,
21         }
22         self.search_engine = SearchEngine(self,
23                                         PolicySpace)
24 ...
25 def run(self):
26     #user code for compilation and execution
27     ...
28     return Result(time=run_result['time'])
29 ...
30 if __name__ == '__main__':
31     opt = SuperOptimizer(statefs, tranfs, actions,
32                         rewards, rl_algs)
33     policy = opt.policy_search(
34         training_benchmark_list, num_of_trials=100)

```

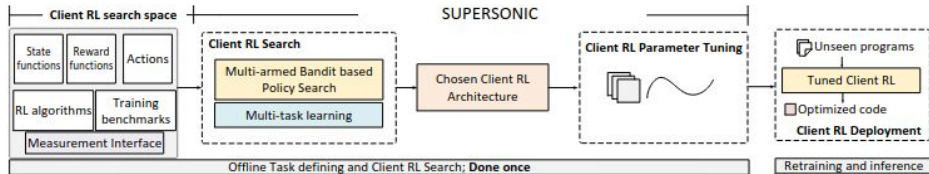


Figure 1. Overview of SUPERSONIC components. This framework enables developers to express the optimization space. It automatically searches for the optimal client RL architecture to be used for inference during deployment.

Table 2. Case studies in our evaluation

Use cases	#Bench.	Competing Methods	Search space
C1: Optimizing image pipeline	10	Halide master [62], auto-schedulers [3], HalideRL [68], OpenTuner	$7^4 * 2^7 \sim 7^4 * 3^3$
C2: Neural network code optimization	5	AutoTVM [19], Chameleon[5], OpenTuner	$5 * 10^3 \sim 20 * 10^4$
C3: Code size reduction	43	CompilerGym [23], OpenTuner	$123^{60} \sim 123^{156}$
C4: Superoptimization	40	STOKE [76], OpenTuner [9]	$9100,000 \sim 916,000,000$

Table 3. Candidate state functions and reward functions

	Case Studies: 1	2	3	4
State func.				
Word2Vec [58]				
Doc2Vec [49]				
CodeBERT [29]				
Manual features (e.g. LLVM IR representation from [40])		✓		✓
Action History				✓
Hash of Action History				✓
Relative measure (e.g. speedup, code size reduction ratio)		✓	✓	✓
Reward func.				
Function output (e.g. bool)				

Table 5. Search overhead required by SUPERSONIC to exceed the performance given the best-performing alternative

Use cases	% of search time (& raw numbers)			% of iterations (& raw numbers)		
	MIN	GeoMean	MAX	MIN	GeoMean	MAX
Case study 1 (21 mins)	3.6	39.4	72.1	3.9	44.0	75.3
Case 1.0	1.0	39.1	74.1	4.2	32.3	68.8
study 2 (1 min)	1.5	31.0	61.1	2.1	29.7	83.5
Case 1.1	1.1	32.8	42.3	1.2	36.5	46.9
study 4 (22 sec)	1.1	32.8	42.3	1.2	36.5	46.9

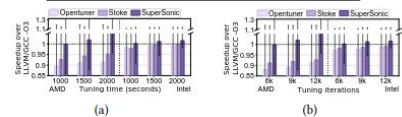


Figure 6. Speedup over LLVM/GCC -O3 for superoptimization under different search time (a) and iterations (b).

Table 1. Example tunable parameters

Algorithms	Parameters
Common param.	Batch size for workers; Train batch size; Mini-batch size; learning rate; Dirichlet noise and epsilon; Puct coefficient; Simulation times; Loss temperature
MCTS	Entropy coefficient; Adam optimizer; step size; GAE estimator parameter; Policy ratio clipping; softmax; Discrete supports; Adam epsilon; Clip gradients
PPO	Behavior Cloning; Pretraining; numbers; Q-Learning loss temperature; Lagrangian threshold; Min Q weight multiplier
DQN	
Q-Learning	

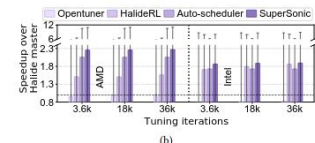
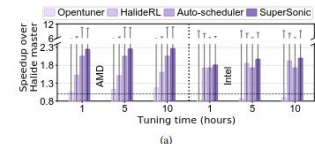
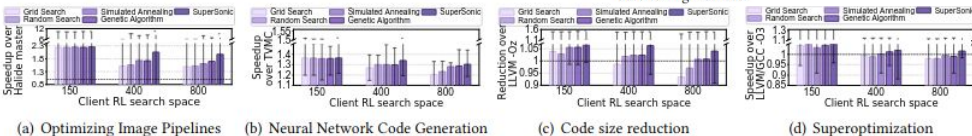


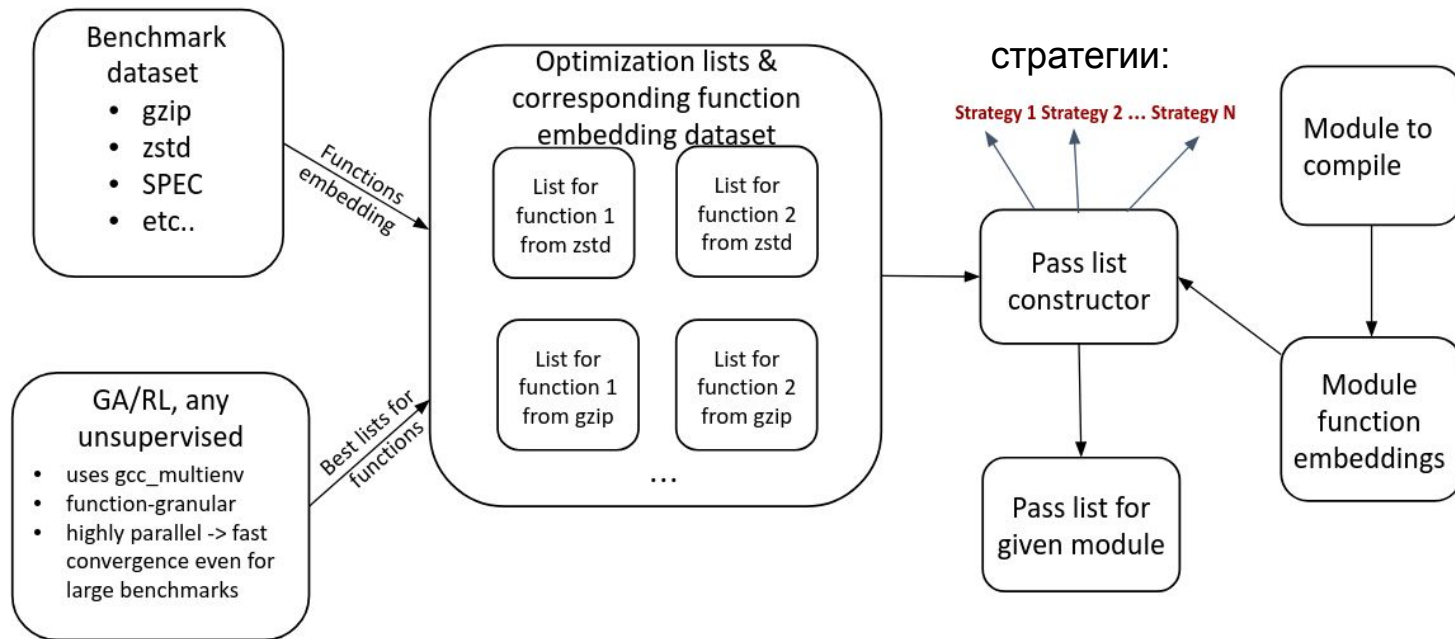
Figure 3. Performance to expert-tuned Halide schedules under different search time (a) and iteration (b) constraints. SUPERSONIC gives the overall best performance than other auto-tuning methods.



- Code features language models are used -- does not take into account data flow and control flow
- Mainly implements Iterative approaches
- I think that this is some overkill for us
- + Many different techniques and methods for finding the optimum are integrated, automated decision-making, which technique to use in a particular case
- + Seems as very powerful and flexible meta-optimization framework for different ML-driven code transformation & optimization tasks

RL/GA combining with supervised learning

RL / GA + supervised learning

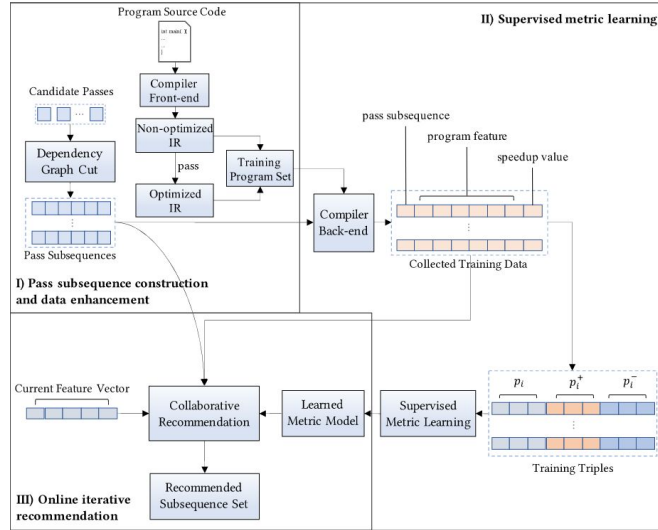


Scenarios:

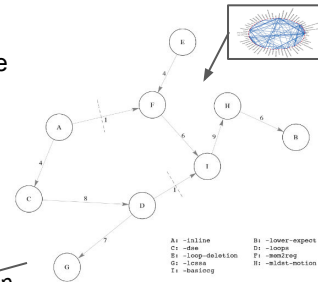
- **Program classification by classes with known best optimizations**
- **Aggregation of passes to greater granularity ones (см. Рис.)**
 - Best optimizations for functions are known
 - Aggregate the explicit sequence for module compilation (by passing this sequence as options to compiler driver / optimizer):
 - Compilation speed is important (but not critical)
 - Goal is to minimize the loss of quality
 - Mainly, controllability and security (trustworthiness) reasons

Iterative Compilation Optimization Based on Metric Learning and Collaborative Filtering (2021) [1]

Using the recommendation system approach: learn the metrics between the program representations in that way that programs with similar best pass sequences be closer in embedding space each to another.



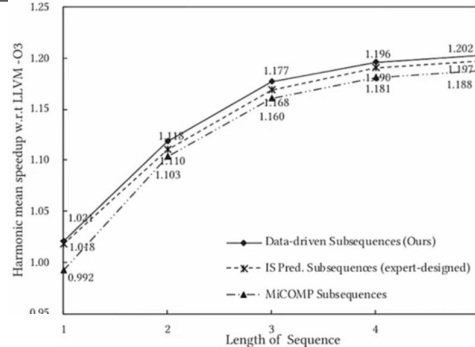
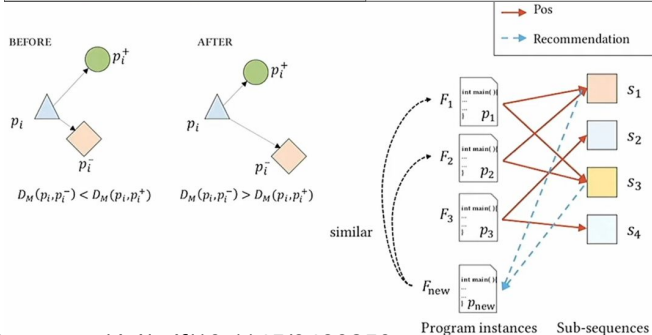
- Consider that pass sequence (A,B) is preferable for program p than (A) and (B) separately, and also than (B,A) , if $f(p, (A,B)) > f(p, (A)) ; f(p, (A,B)) > f(p, (B)) ; f(p, (A,B)) > f(p, (B,A))$, where f is the metric of performance gain.
- Such passes, that A improves effectiveness of pass B , so B depends on A , are called “collaborative interactions” between A and B .
- As collaborative interactions are oriented 2-ary relations on passes, the full set of explored pass relations can be represented as oriented dependency graph -- ODG.



Collaborative filtering:

Recommend on the inference such subsequences that best optimize programs “similar to the current program” in some metric.

[6] : learn metric such that the distance between programs for which the same pass sequences fit is small, and for not same it is large



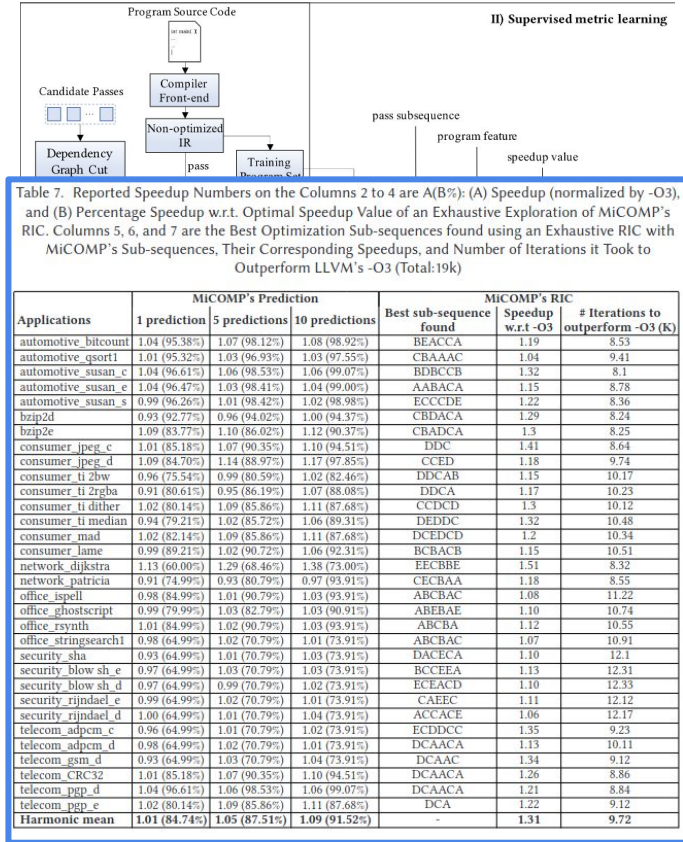
The article is based on more earlier article [1] with improved Approach for subsequences construction. The [2] based on graph-structured agglomeration, when collaborative filtering with metric learning in [1] gives better results.

[1] Hongzhi Liu, Jie Luo, Ying Li, and Zhonghai Wu. 2021. Iterative Compilation Optimization Based on Metric Learning and Collaborative Filtering. ACM Trans. Arch. Code Optim. 19, 1, Article 2 (December 2021), 25 pages. <https://doi.org/10.1145/3480250>

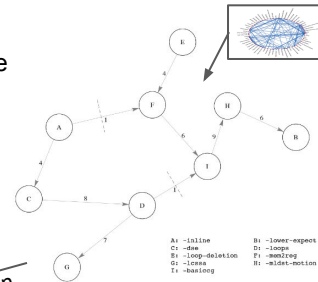
[2] Amir H. Ashouri, Andrea Bignoli, Gianluca Palermo, Cristina Silvano, Sameer Kulkarni, and John Cavazos. 2017. MiCOMP: Mitigating the Compiler Phase-Ordering Problem Using Optimization Sub-Sequences and Machine Learning. ACM Trans. Archit. Code Optim. 14, 3, Article 29 (September 2017), 28 pages. <https://doi.org/10.1145/3124452>

Iterative Compilation Optimization Based on Metric Learning and Collaborative Filtering (2021) [1]

Using the recommendation system approach: learn the metrics between the program representations in that way that programs with similar best pass sequences be closer in embedding space each to another.

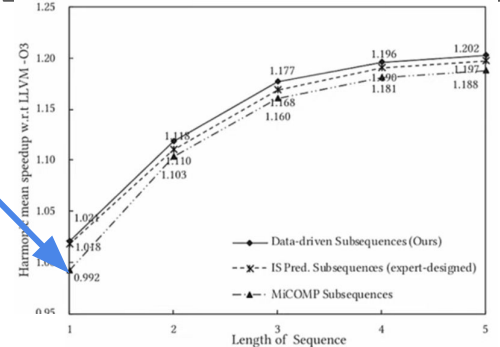


- Consider that pass sequence (A,B) is preferable for program p than (A) and (B) separately, and also than (B,A), if $f(p, (A,B)) > f(p, (A)) ; f(p, (A,B)) > f(p, (B)) ; f(p, (A,B)) > f(p, (B,A))$, where f is the metric of performance gain.
- Such passes, that A improves effectivity of pass B, so B depends on A, are called "collaborative interactions" between A and B.
- As collaborative interactions are oriented 2-ary relations on passes, the full set of explored pass relations can be represented as oriented dependency graph -- ODG.



Collaborative filtering: Recommend on the inference such subsequences that best optimize programs "similar to the current program" in some metric.

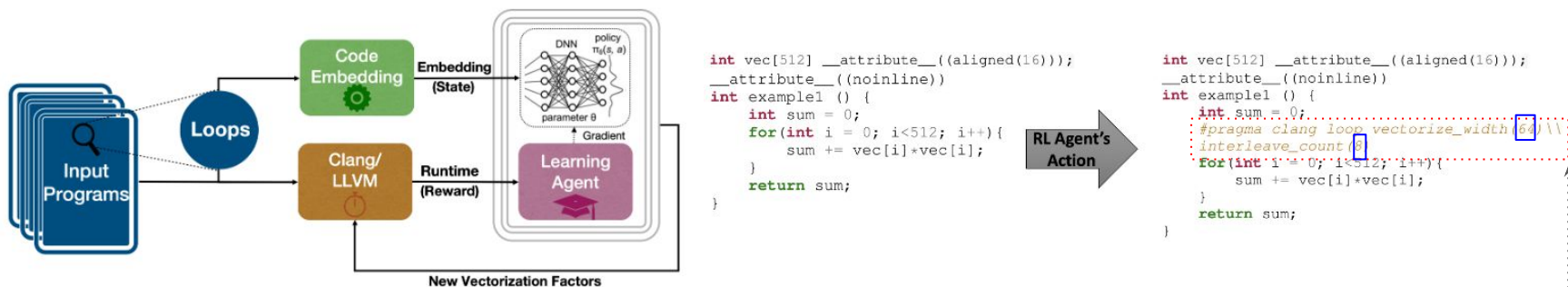
[6] : learn metric such that the distance between programs for which the same pass sequences fit is small, and for not same it is large



The article is based on more earlier article [1] with improved approach for subsequences construction. The [2] based on graph-structured agglomeration, when collaborative filtering with metric learning in [1] gives better results. [1] Hongzhi Liu, Jie Luo, Ying Li, and Zhonghai Wu. 2021. Iterative Compilation Optimization Based on Metric Learning and Collaborative Filtering. ACM Trans. Arch. Code Optim. 19, 1, Article 2 (December 2021), 25 pages. <https://doi.org/10.1145/3480250> [2] Amir H. Ashouri, Andrea Bignoli, Gianluca Palermo, Cristina Silvano, Sameer Kulkarni, and John Cavazos. 2017. MiCOMP: Mitigating the Compiler Phase-Ordering Problem Using Optimization Sub-Sequences and Machine Learning. ACM Trans. Archit. Code Optim. 14, 3, Article 29 (September 2017), 28 pages. <https://doi.org/10.1145/3124452>

Pass parameters tuning,
Arch-dependent optimizations

Pass parameters tuning / Arch-dependent optimizations



NeuroVectorizer** RL agent's step example

- Mostly unsupervised
- Most simple characterization -- максимально простые
 - Expertly selected features
 - Sometimes code2vec (NLP-based)
- Usually tune a few optimizations (oftenly, only one, but much accurately)
- LLVM-based
- Examples:
 - MLGO (2021)*: <https://github.com/google/ml-compiler-opt>
 - NeuroVectorizer (2020)**: <https://github.com/intel/neuro-vectorizer>

* <https://blog.research.google/2022/07/mlgo-machine-learning-framework-for.html>

Summary:

- There are a number of not only researched, but also tested in practice solutions for iterative pass order auto-tuning.
- Most solutions look for passes for the entire translation unit
- Multi-step solutions look preferable, in particular, the superiority indicated by the authors of the articles even over expertly selected sequences, while reducing the search time and overhead for iterative measurement. Maybe it deserves consideration ...
- Splitting into subsequences in a multi-step approach is also an open question
- None of the solutions take into account the multi-criteria choice
- None of the phase-ordering/choosing approaches tune the passes parametrization (for ex. 'UF' for loop-unroll)
- Most of solutions use a quite inaccurate model for representing (characterizing) programs
 - Ignoring information about data flows
 - Ignoring the control flow
 - Operating:
 - By fixed set of features (for example, the number of instructions of a given type, the size of the BB, etc. - for example, in [1] there are 56 of them, but aggregated as a sum for each of the components of the feature vector)
 - Or by embedding built on mechanisms from the field of natural language processing (word2vec, code2vec, CodeBERT, etc)
 - LLMs are quite new to use it

Review conclusion:

- The most proven solution in practice for now is AutoPhase [1].
 - The used model of the program features [1] is primitive -- IR2Vec (or alternative, for ex. PrograML[11]) required [4], taking into account data and control flows.
 - It is important for convergency speed-up
 - Improvements are also required in multi-step mode -- selection of subsequences from ODG, and different policies for its implementation [2,3,6].
- None of the solutions take into account the multi-criteria choice, only few solutions are constrained. Some scenarios requires combination with supervised learning
- None of the solutions takes into account the parametrization of passes (Except MLGO and NeuroVectorizers, which focused only on few passes)
- Additionally, the issue of optimizing individual functions (partial compilation) can be considered

Related work

1. Q. Huang, et al., "AutoPhase: Compiler Phase-Ordering for HLS with Deep Reinforcement Learning," in 2019 IEEE 27th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM), San Diego, CA, USA, 2019 pp. 308-308.
2. POSET-RL: <https://arxiv.org/abs/2208.04238>
3. [Wang, H, Tang, Z, Zhang, C et al]. (4 more authors) (2022) Automating Reinforcement Learning Architecture Design for Code Optimization. In: CC 2022: Proceedings of the 31st ACM SIGPLAN International Conference on Compiler Construction. The 31st ACM SIGPLAN International Conference on Compiler Construction (CC '22), 02-03 Apr 2022.
4. S. VenkataKeerthy, Rohit Aggarwal, Shalini Jain, Maunendra Sankar Desarkar, Ramakrishna Upadrasta, and Y. N. Srikant. 2020. IR2VEC: LLVM IR Based Scalable Program Embeddings. ACM Trans. Archit. Code Optim. 17, 4, Article 32 (December 2020), 27 pages.
5. Chris Cummins, Bram Wasti, Jiadong Guo, Brandon Cui, Jason Ansel, Sahir Gomez, Somya Jain, Jia Liu, Olivier Teytaud, Benoit Steiner, Yuandong Tian, and Hugh Leather. 2022. CompilerGym: robust, performant compiler optimization environments for AI research. In Proceedings of the 20th IEEE/ACM International Symposium on Code Generation and Optimization (CGO '22). IEEE Press, 92–105.
6. Liu et al. (2021) Hongzhi Liu, Jie Luo, Ying Li, and Zhonghai Wu. 2021. Iterative Compilation Optimization Based on Metric Learning and Collaborative Filtering. ACM Trans. Arch. Code Optim. 19, 1, Article 2 (December 2021), 25 pages. <https://doi.org/10.1145/3480250>
7. Amir H. Ashouri, Andrea Bignoli, Gianluca Palermo, Cristina Silvano, Sameer Kulkarni, and John Cavazos. 2017. MiCOMP: Mitigating the Compiler Phase-Ordering Problem Using Optimization Sub-Sequences and Machine Learning. ACM Trans. Archit. Code Optim. 14, 3, Article 29 (September 2017), 28 pages. <https://doi.org/10.1145/3124452>
8. Zavodskikh R. K., Efanov N. N., Tomashev D. D. (2022). Using the LLVM Framework for static performance prediction with embedding of intermediate representation. Proceedings of Moscow Institute of Physics and Technology. Vol. 14, 3 (55) pp 34-45.
9. Tağtekin, B., Höke, B., Sezer, M. K., & Öztürk, M. U. (2021, August). FOGA: Flag Optimization with Genetic Algorithm. In *2021 International Conference on INnovations in Intelligent SysTems and Applications (INI/STA)* (pp. 1-6). IEEE.
10. F. Agakov, E. Bonilla, J. Cavazos, B. Franke, G. Fursin, M. F. O'Boyle, J. Thomson, M. Toussaint, and C. K. Williams, "Using machine learning to focus iterative optimization," in International Symposium on Code Generation and Optimization (CGO'06). IEEE, 2006, pp. 11–pp.
11. PROGRAML: A Graph-based Program Representation for Data Flow Analysis and Compiler Optimizations (2021) <https://chriscummins.cc/pub/2021-icml.pdf>