

Lecture Notes: Program Synthesis

17-355/17-665/17-819: Program Analysis (Spring 2019)

Jonathan Aldrich*

aldrich@cs.cmu.edu

Note: A complete, if lengthy, resource on inductive program synthesis is the book “Program Synthesis” by Gulwani *et. al* [8]. You need not read the whole thing; I encourage you to investigate the portions of interest to you, and skim as appropriate. Many references in this document are drawn from there; if you are interested, it contains many more.

1 Program Synthesis Overview

The problem of program synthesis can be expressed as follows:

$$\exists P . \forall x . \varphi(x, P(x))$$

In the setting of *constructive* logic, proving the validity of a formula that begins with an existential involves coming up with a *witness*, or concrete example, that can be plugged into the rest of the formula to demonstrate that it is true. In this case, the witness is a program P that satisfies some specification φ on all inputs. We take a liberal view of P in discussing synthesis, as a wide variety of artifact types have been successfully synthesized (anything that reads inputs or produces outputs). Beyond (relatively small) program snippets of the expected variety, this includes protocols, interpreters, classifiers, compression algorithms or implementations, scheduling policies, and cache coherence protocols for multicore processors. The specification φ is an expression of the user intent, and may be expressed in one of several ways: a formula, a reference implementation, input/output pairs, traces, demonstrations, or a syntactic *sketch*, among other options.

Program synthesis can thus be considered along three dimensions:

(1) Expressing user intent. User intent (or φ in the above) can be expressed in a number of ways, including logical specifications, input/output examples [4] (often with some kind of user- or synthesizer-driven interaction), traces, natural language [3, 7, 13], or full- or partial programs [20]. In this latter category lies reference implementations, such as executable specifications (which give the desired output for a given input) or declarative specifications (which check whether a given input/output pair is correct). Some synthesis techniques allow for multi-modal specifications, including pre- and post- conditions, safety assertions at arbitrary program points, or partial program templates.

Such specifications can constrain two aspects of the synthesis problem:

- **Observable behavior**, such as an input/output relation, a full executable specification or safety property. This specifies *what* a program should compute.

*These notes are adapted from an earlier version written by Claire Le Goues

- **Structural properties**, or internal computation steps. These are often expressed as a sketch or template, but can be further constrained by assertions over the number or variety of operations in a synthesized programs (or number of iterations, number of cache misses, etc, depending on the synthesis problem in question). Indeed, one of the key principles behind the scaling of many modern synthesis techniques lie in the way they syntactically restrict the space of possible programs, often via a sketch, grammar, or DSL.

Note that basically all of the above types of specifications can be translated to constraints in some form or another. Techniques that operate over multiple types of specifications can overcome various challenges that come up over the course of an arbitrary synthesis problem. Different specification types are more suitable for some types of problems than others. In addition, trace- or sketch-based specifications can allow a synthesizer to decompose a synthesis problems into intermediate program points.

Question: how many ways can we specify a sorting algorithm?

(2) Search space of possible programs. The search space naturally includes programs, often constructed of subsets of normal programming languages. This can include a predefined set of considered operators or control structures, defined as grammars. However, other spaces are considered for various synthesis problems, like logics of various kinds, which can be useful for, e.g., synthesizing graph/tree algorithms.

(3) Search technique. At a high level, there are two general approaches to logical synthesis:

- Deductive (or classic) synthesis (e.g., [15]), which maps a high-level (e.g. logical) specification to an executable implementation, classically using a theorem prover. Such approaches are efficient and provably correct: thanks to the semantics-preserving rules, only correct programs are explored. However, they require complete specifications and sufficient axiomatization of the domain. These approaches are classically applied to e.g., controller synthesis.
- Inductive (sometimes called syntax-guided) synthesis, which takes a partial (and often multi-modal) specification and constructs a program that satisfies it. These techniques are more flexible in their specification requirements and require no axioms, but often at the cost of lower efficiency and weaker bounded guarantees on the optimality of synthesized code.

Deductive synthesis shares quite a bit in common, conceptually, with compilation: rewriting a specification according to various rules to achieve a new program in at a different level of representation. However, deductive synthesis approaches assume a complete formal specification of the desired user intent was provided. In many cases, this can be as complicated as writing the program itself.

This has motivated new inductive synthesis approaches, towards which considerable modern research energy has been dedicated. This category of techniques lends itself to a wide variety of search strategies, including brute-force or enumerative [1] (you might be surprised!), probabilistic inference/belief propagation [6], or genetic programming [12]. Alternatively, techniques based on logical reasoning delegate the search problem to a constraint solver. We will spend more time on this set of techniques.

2 Deductive Synthesis

We will very briefly overview Denali [11], a prototypical deductive synthesis technique for *superoptimization*.¹ Denali seeks to generate short sequences of provably optimal loop-free machine instructions, for use primarily in compilation. While compilers generate reasonably good code, there are cases in which we would instead prefer provably optimal code. Generating such code is the task of a superoptimizer (so-called because the title of optimization “has been given to a field that does not aspire to optimize but only to improve”). Early approaches for superoptimization attempted to enumerate via brute force (in order of increasing length) efficient sequences of instructions, with correctness checked by hand and against a set of test cases. This correctness criterion is challenging to confirm, however, and does not necessarily result in optimality.

Joshi et al. propose an approach to superoptimization based on theorem proving. The “obvious” approach (which they do *not* take) would be to, given a desired program fragment P , express in formal logic “no program of the target architecture computes P in at most N cycles.” However, this obvious approach is very difficult to manage with a theorem prover, because it must be expressed using nested quantifiers.

Instead, they propose a process based on the idea that for sufficiently simple programs, equivalence between a desired P and some alternative implementation M for all inputs is essentially the universal validity of an equality between two vectors of terms (the one M computes, and the terms specified by P in the computation). This type of equivalence can be proved by matching, which is a well understood technique in theorem proving.

To do this, their technique, named Denali, takes as input a program P written in a DSL for the associated target architecture. It then constructs an *E-graph* using the specified desired program P as input. An E-graph is a term DAG corresponding to the expression to be synthesized, augmented with an equivalence relation on the nodes of the DAG. Two nodes are equivalent if the terms they represent are identical in value. Denali then uses a theorem prover, along with two sets of axioms (encoding *instruction semantics*—an interpreter for the target language, effectively—and algebraic properties—memory modeling, mostly), to search the e-graph for the most efficient way to compute the expression.

3 Inductive Synthesis

Inductive synthesis uses inductive reasoning to construct programs in response to partial specifications. The program is synthesized via a symbolic interpretation of a space of candidates, rather than by deriving the candidate directly. So, to synthesize such a program, we basically only require an interpreter, rather than a sufficient set of derivation axioms. Inductive synthesis is applicable to a variety of problem types, such as string transformation (FlashFill) [5], data extraction/processing/wrangling [4, 19], layout transformation of tables or tree-shaped structures [21], graphics (constructing structured, repetitive drawings) [9, 2], program repair [16, 14] (spoiler alert!), superoptimization [11], and efficient synchronization, among others.

Inductive synthesis consists of several family of approaches; we will overview several prominent examples, without claiming to be complete.

¹This explanation is further illustrated using the associated lecture slides.

3.1 SKETCH, CEGIS, and SyGuS

SKETCH is a well-known synthesis system that allows programs to provide partial programs (a sketch) that expresses the high-level structure of the intended implementation but leaves holes for low-level implementation details. The synthesizer fills these holes from a finite set of choices, using an approach now known as Counterexample-guided Inductive Synthesis (CEGIS) [20, 18]. This well-known synthesis architecture divides the problem into *search* and *verification* components, and uses the output from the latter to refine the specification given to the former.

We have a diagram to illustrate on slides.

Syntax-Guided Synthesis (or *SyGuS*) formalizes the problem of program synthesis where specification is supplemented with a syntactic template. This defines a search space of possible programs that the synthesizer effectively traverses. Many search strategies exist; two especially well-known strategies are *enumerative search* (which can be remarkably effective, though rarely scales), and *deductive* or *top down* search, which recursively reduces the problem into simpler sub-problems.

3.2 Oracle-guided synthesis

Templates or sketches are often helpful and easy to write. However, they are not always available. Beyond search or enumeration, constraint-based approaches translate a program's specification into a constraint system that is provided to a solver. This can be especially effective if combined with an outer CEGIS loop that provides oracles.

This kind of synthesis can be effective when the properties we care about are relatively easy to verify. For example, imagine we wanted to find a maximum number m in a list l .

Turn to the handout, which asks you to specify this as a synthesis problem...

Note that instead of proving that a program satisfies a given formula, we can instead disprove its negation, which is:

$$\exists l, m : (P_{max}(l) = m) \wedge (m \notin l \vee \exists x \in l : m < x)$$

If the above is satisfiable, a solver will give us a counterexample, which we can use to strengthen the specification—so that next time the synthesis engine will give us a program that excludes this counterexample. We can make this counterexample more useful by asking the solver not just to provide us with an input that produces an error, but also to provide the corresponding correct output m^* :

$$\exists l, m^* : (P_{max}(l) \neq m^*) \wedge (m^* \in l) \wedge (\forall x \in l : m^* \geq x)$$

This is a much stronger constraint than the original counterexample, as it says what the program should output in this case rather than one example of something it should not output. Thus we now have an additional test case for the next round of synthesis. This counterexample-guided synthesis approach was originally introduced for SKETCH, and was generalized to oracle-guided inductive synthesis by Jha and Seshia. Different oracles have been developed for this type of synthesis. We will discuss component-based oracle-guided program synthesis in detail, which illustrates the use of distinguishing oracles.

4 Oracle-guided Component-based Program Synthesis

Problem statement and intuition.² Given a set of input-output pairs $\langle \alpha_0, \beta_0 \rangle \dots \langle \alpha_n, \beta_n \rangle$ and N components f_1, \dots, f_n , the goal is to synthesize a function f out of the components such that $\forall \alpha_i. f(\alpha_i)$ produces β_i . We achieve this by constructing and solving a set of constraints over f , passing those constraints to an SMT solver, and using a returned satisfying model to reconstruct f . In this approach, the synthesized function will have the following form:

```

0           $z_0 := \text{input}^0$ 
1           $z_1 := \text{input}^1$ 
...
 $m$         $z_m := \text{input}^m$ 
 $m + 1$     $z_{m+1} := f_?(z_?, \dots, z_?)$ 
 $m + 2$     $z_{m+2} := f_?(z_?, \dots, z_?)$ 
...
 $m + n$     $z_{m+n} := f_?(z_?, \dots, z_?)$ 
 $m + n + 1$  return  $z_?$ 

```

The thing we have to do is fill in the ? indexes in the program above. These indexes essentially define the order in which functions are invoked and what arguments they are invoked with. We will assume that each component is used once, without loss of generality, since we can duplicate the components.

Definitions. We will set up the problem for the solver using two sets of variables. One set represents the input values passed to each component, and the output value that component produces, when the program is run for a given test case. We use $\vec{\chi}_i$ to denote the vector of input values passed to component i and r_i to denote the result value computed by that component. So if we have a single component (numbered 1) that adds two numbers, the input values $\vec{\chi}_1$ might be (1,3) for a given test case and the output r_1 in that case would be 4. We use Q to denote the set of all variables representing inputs and R to denote the set of all variables representing outputs:

$$\begin{aligned}
 Q &:= \bigcup_{i=1}^N \vec{\chi}_i \\
 R &:= \bigcup_{i=1}^N r_i
 \end{aligned}$$

We also define the overall program's inputs to be the vector \vec{Y} and the program's output to be r .

The other set of variables determines the location of each component, as well as the locations at which each of its inputs were defined. We call these *location variables*. For each variable x , we define a location variable l_x , which denotes where x is defined. Thus l_{r_i} is the location variable for the result of component i and \vec{l}_{χ_i} is the vector of location variables for the inputs of component i . So if we have $l_{r_3} = 5$ and \vec{l}_{χ_3} is (2,4), then we will invoke component #3 at line 5, and we will pass variables z_2 and z_4 to it. L is the set of all location variables:

$$L := \{l_x | x \in Q \cup R \cup \vec{Y} \cup r\}$$

We will have two sets of constraints: one to ensure the program is *well-formed*, and the other that ensures the program encodes the desired *functionality*.

²These notes are inspired by Section III.B of Nguyen *et al.*, ICSE 2013 [17] ...which provides a really beautifully clear exposition of the work that originally proposed this type of synthesis in Jha *et al.*, ICSE 2010 [10].

Well-formedness. ψ_{wfp} denotes the well-formedness constraint. Let $M = |\vec{Y}| + N$, where N is the number of available components:

$$\psi_{wfp}(L, Q, R) \stackrel{def}{=} \bigwedge_{x \in Q} (0 \leq l_x < M) \wedge \bigwedge_{x \in R} (|\vec{Y}| \leq l_x < M) \wedge \psi_{cons}(L, R) \wedge \psi_{acyc}(L, Q, R)$$

The first line of that definition says that input locations are in the range 0 to M , while component output locations are all defined after program inputs are declared. ψ_{cons} and ψ_{acyc} dictate that there is only one component in each line and that the inputs of each component are defined before they are used, respectively:

$$\begin{aligned} \psi_{cons}(L, R) &\stackrel{def}{=} \bigwedge_{x, y \in R, x \neq y} (l_x \neq l_y) \\ \psi_{acyc}(L, Q, R) &\stackrel{def}{=} \bigwedge_{i=1}^N \bigwedge_{x \in \vec{\chi}_i, y \equiv r_i} l_x < l_y \end{aligned}$$

Functionality. ϕ_{func} denotes the functionality constraint that guarantees that the solution f satisfies the given input-output pairs:

$$\begin{aligned} \phi_{func}(L, \alpha, \beta) &\stackrel{def}{=} \psi_{conn}(L, \vec{Y}, r, Q, R) \wedge \phi_{lib}(Q, R) \wedge (\alpha = \vec{Y}) \wedge (\beta = r) \\ \psi_{conn}(L, \vec{Y}, r, Q, R) &\stackrel{def}{=} \bigwedge_{x, y \in Q \cup R \cup \vec{Y} \cup \{r\}} (l_x = l_y \Rightarrow x = y) \\ \phi_{lib}(Q, R) &\stackrel{def}{=} \left(\bigwedge_{i=1}^N \phi_i(\vec{\chi}_i, r_i) \right) \end{aligned}$$

ψ_{conn} encodes the meaning of the location variables: If two locations are equal, then the values of the variables defined at those locations are also equal. ϕ_{lib} encodes the semantics of the provided basic components, with ϕ_i representing the specification of component f_i . The rest of ϕ_{func} encodes that if the input to the synthesized function is α , the output must be β .

Almost done! ϕ_{func} provides constraints over a single input-output pair α_i, β_i , we still need to generalize it over all n provided pairs $\{ \langle \alpha_i, \beta_i \rangle \mid 1 \leq i \leq n \}$:

$$\theta \stackrel{def}{=} \left(\bigwedge_{i=1}^n \phi_{func}(L, \alpha_i, \beta_i) \right) \wedge \psi_{wfp}(L, Q, R)$$

θ collects up all the previous constraints, and says that the synthesized function f should satisfy all input-output pairs and the function has to be well formed.

LVal2Prog. The only real unknowns in all of θ are the values for the location variables L . So, the solver that provides a satisfying assignment to θ is basically giving a valuation of L that we then turn into a constructed program as follows:

Given a valuation of L , Lval2Prog(L) converts it to a program as follows: The i^{th} line of the program is $z_i = f_j(z_{\sigma_1}, \dots, r_{\sigma_\eta})$ when $l_{r_j} = i$ and $\bigwedge_{k=1}^{\eta} (l_{\chi_j^k} = \sigma_k)$, where η is the number of inputs for component f_j and χ_j^k denotes the k^{th} input parameter of component f_j . The program output is produced in line l_r .

Example. Assume we only have one component, $+$. $+$ has two inputs: χ_+^1 and χ_+^2 . The output variable is r_+ . Further assume that the desired program f has one input Y_0 (which we call input^0 in the actual program text) and one output r . Given a mapping for location variables of: $\{l_{r_+} \mapsto 1, l_{\chi_+^1} \mapsto 0, l_{\chi_+^2} \mapsto 0, l_r \mapsto 1, l_Y \mapsto 0\}$, then the program looks like:

```

0   $z_0 := \text{input}^0$ 
1   $z_1 := z_0 + z_0$ 
2  return  $z_1$ 

```

This occurs because the location of the variables used as input to $+$ are both on the same line (0), which is also the same line as the input to the program (0). l_r , the return variable of the program, is defined on line 1, which is also where the output of the $+$ component is located. (l_{r_+}). We added the return on line 2 as syntactic sugar.

References

- [1] R. Alur, R. Bodík, E. Dallal, D. Fisman, P. Garg, G. Juniwal, H. Kress-Gazit, P. Madhusudan, M. M. K. Martin, M. Raghothaman, S. Saha, S. A. Seshia, R. Singh, A. Solar-Lezama, E. Torlak, and A. Udupa. Syntax-guided synthesis. In M. Irlbeck, D. A. Peled, and A. Pretschner, editors, *Dependable Software Systems Engineering*, volume 40 of *NATO Science for Peace and Security Series, D: Information and Communication Security*, pages 1–25. IOS Press, 2015.
- [2] R. Chugh, B. Hempel, M. Spradlin, and J. Albers. Programmatic and direct manipulation, together at last. *SIGPLAN Not.*, 51(6):341–354, June 2016.
- [3] A. Desai, S. Gulwani, V. Hingorani, N. Jain, A. Karkare, M. Marron, S. R., and S. Roy. Program synthesis using natural language. In *Proceedings of the 38th International Conference on Software Engineering, ICSE ’16*, pages 345–356, New York, NY, USA, 2016. ACM.
- [4] S. Gulwani. Programming by examples: Applications, algorithms, and ambiguity resolution. In *Proceedings of the 8th International Joint Conference on Automated Reasoning - Volume 9706*, pages 9–14, New York, NY, USA, 2016. Springer-Verlag New York, Inc.
- [5] S. Gulwani, W. R. Harris, and R. Singh. Spreadsheet data manipulation using examples. *Commun. ACM*, 55(8):97–105, Aug. 2012.
- [6] S. Gulwani and N. Jojic. Program verification as probabilistic inference. In *Proceedings of the 34th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL ’07*, pages 277–289, New York, NY, USA, 2007. ACM.
- [7] S. Gulwani and M. Marron. Nlyze: Interactive programming by natural language for spreadsheet data analysis and manipulation. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data, SIGMOD ’14*, pages 803–814, New York, NY, USA, 2014. ACM.
- [8] S. Gulwani, O. Polozov, and R. Singh. Program synthesis. *Foundations and Trends in Programming Languages*, 4(1-2):1–119, 2017.
- [9] B. Hempel and R. Chugh. Semi-automated svg programming via direct manipulation. In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology, UIST ’16*, pages 379–390, New York, NY, USA, 2016. ACM.

- [10] S. Jha, S. Gulwani, S. A. Seshia, and A. Tiwari. Oracle-guided component-based program synthesis. In *Proceedings of the 32Nd ACM/IEEE International Conference on Software Engineering - Volume 1, ICSE '10*, pages 215–224, New York, NY, USA, 2010. ACM.
- [11] R. Joshi, G. Nelson, and K. Randall. Denali: A goal-directed superoptimizer. *SIGPLAN Not.*, 37(5):304–314, May 2002.
- [12] G. Katz and D. Peled. Genetic programming and model checking: Synthesizing new mutual exclusion algorithms. In *Proceedings of the 6th International Symposium on Automated Technology for Verification and Analysis, ATVA '08*, pages 33–47, Berlin, Heidelberg, 2008. Springer-Verlag.
- [13] V. Le, S. Gulwani, and Z. Su. Smartsynth: Synthesizing smartphone automation scripts from natural language. In *Proceeding of the 11th Annual International Conference on Mobile Systems, Applications, and Services, MobiSys '13*, pages 193–206, New York, NY, USA, 2013. ACM.
- [14] C. Le Goues, T. Nguyen, S. Forrest, and W. Weimer. GenProg: A generic method for automated software repair. *IEEE Transactions on Software Engineering*, 38(1):54–72, 2012.
- [15] Z. Manna and R. J. Waldinger. Toward automatic program synthesis. *Commun. ACM*, 14(3):151–165, Mar. 1971.
- [16] S. Mechtaev, J. Yi, and A. Roychoudhury. Angelix: Scalable Multiline Program Patch Synthesis via Symbolic Analysis. In *International Conference on Software Engineering, ICSE '16*, pages 691–701, 2016.
- [17] H. D. T. Nguyen, D. Qi, A. Roychoudhury, and S. Chandra. Semfix: Program repair via semantic analysis. In *Proceedings of the 2013 International Conference on Software Engineering, ICSE '13*, pages 772–781, Piscataway, NJ, USA, 2013. IEEE Press.
- [18] O. Polozov and S. Gulwani. Flashmeta: A framework for inductive program synthesis. *SIGPLAN Not.*, 50(10):107–126, Oct. 2015.
- [19] R. Singh and S. Gulwani. Transforming spreadsheet data types using examples. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '16*, pages 343–356, New York, NY, USA, 2016. ACM.
- [20] A. Solar-Lezama. *Program Synthesis by Sketching*. PhD thesis, Berkeley, CA, USA, 2008. AAI3353225.
- [21] N. Yaghmazadeh, C. Klinger, I. Dillig, and S. Chaudhuri. Synthesizing transformations on hierarchically structured data. *SIGPLAN Not.*, 51(6):508–521, June 2016.