



**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

MASTER THESIS

Petr Onderka

**System for extensions
of the C# language**

Department of Software Engineering

Supervisor of the master thesis: RNDr. Filip Zavoral, Ph.D.

Study programme: Computer Science

Study branch: Software and Data Engineering

Prague 2018

I declare that I carried out this master thesis independently, and only with the cited sources, literature and other professional sources.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In date

signature of the author

Dedication.

Title: System for extensions of the C# language

Author: Petr Onderka

Department: Department of Software Engineering

Supervisor: RNDr. Filip Zavoral, Ph.D., Department of Software Engineering

Abstract: Abstract.

Keywords: key words

Contents

Introduction	3
1 Background	4
1.1 Example	4
1.2 Manipulating C# source code	5
1.2.1 T4	5
1.2.2 CodeDOM	5
1.2.3 Roslyn	6
1.3 Manipulating IL	11
1.3.1 System.Reflection.Emit	11
1.3.2 Mono.Cecil	11
1.4 Other approaches	11
1.4.1 Expression trees	11
1.5 Metaprogramming systems	13
1.5.1 PostSharp	13
1.5.2 Fody	14
1.5.3 F# type providers	14
2 Analysis	15
2.1 Representing code	15
2.2 The system	18
3 Design	21
3.1 CSharpE.Syntax	21
3.1.1 General principles	21
3.1.2 Projects	22
3.1.3 Source files	22
3.1.4 Types	23
3.1.5 Members	23
3.1.6 Statements	24
3.1.7 Expressions	24
3.1.8 References	24
3.2 CSharpE.Transform	24
4 Implementation	25
4.1 Syntax trees	25
4.2 Transformations	25
4.3 IDE integration	25
4.4 Example extensions	25
5 Comparison with existing tools	26
5.1 Reading and writing code	26
5.2 Transforming code	26
6 Future work	27

Conclusion	28
Bibliography	30
List of Figures	31
List of Tables	32
List of Abbreviations	33
A Attachments	34
A.1 First Attachment	34

Introduction

Extensibility is an important feature of a programming language and its associated programming environment, because it allows adding new capabilities to the language. This way, a programmer can mold the language to their specific needs.

A common way to extend many languages is through libraries. Libraries let programmers use code written by someone else, which can be powerful, especially when combined with extensibility features built into programming languages, such as virtual functions or lambdas.

But libraries are restricted to the features offered by the used language, which can limit their usefulness.

This work describes a system for extending the C# language beyond what can be accomplished with libraries by user-provided extensions, which perform transformations of C# source code.

These extensions, themselves written in C#, should be easy to create, when compared with existing similar systems, and they should be efficient enough to be usable with code completion in a code editor or an integrated development environment (IDE), such as Microsoft Visual Studio.

To achieve this, the system is composed of two primary parts: the Syntax tree API (Application programming interface) and the Transformations API.

The Syntax tree API is used to represent the original C# source code, examine it, and modify it. The primary goal of this API is to be easy to use and abstract syntax trees fit that requirement well.

The Transformation API enhances the Syntax tree API by adding methods that split source code transformation into smaller parts. The inputs and outputs of each part are then tracked, which means that, after an initial full execution of the transformation, only the parts of the transformation whose inputs changed have to be re-executed. This is done to improve the performance of the system, especially when run from an IDE.

1. Background

The .NET ecosystem is composed of programming languages (including C#, F# and Visual Basic .NET), [1] .NET implementations (including .NET Framework, .NET Core and Mono), [2] class libraries, commonly distributed through the NuGet package manager, and tooling, including command-line tools and tools integrated into code editors and IDEs.

What unifies all these components is the Common Language Infrastructure (CLI), [3] which specifies binary file format for “assemblies”. These contain compiled .NET code in the form of Intermediate Language (IL) and also metadata associated with this code.

The C# language [4] is an object-oriented programming language which is part of the .NET ecosystem. The C# compiler, code named “Roslyn”, [5] compiles C# source code into a .NET assembly. The compiler can also be used as a class library, which exposes types for programmatically manipulating C# source code.

An assembly, produced by the C# compiler or in some other way, can be executed on a .NET implementation. Each .NET implementation contains a runtime, which is responsible for executing code, and a base class library, which contains basic types used by .NET programs.

Runtimes of .NET implementations are usually using a just-in-time (JIT) compiler, which converts the IL for each method into machine code specific for the current instruction set just before executing that method for the first time.

In the .NET ecosystem, class libraries, which are just .NET assemblies, are commonly distributed thorough the NuGet package manager, [6] because it makes using those libraries easier. And while NuGet is primarily used for regular libraries, which are directly used by programmers in their source code, it can also be used for various kinds of special libraries, such as add-ins for general metaprogramming systems like Fody (more on these in section 1.5), or Roslyn analyzers for detecting issues with source code.

The C# language contains some basic extensibility features itself, namely virtual methods and delegates. But for more advanced use cases, it is necessary ot manipulate code in some form and the .NET ecosystem has various approaches to achieve that, including those that manipulate C# code, those that manipulate IL and those that use a custom model for representing code. Some of these approaches will be described in following sections.

1.1 Example

To demonstrate various code generation approaches, a running example of generating a simple entity class with a set of properties and having `IEquatable<T>` as an implemented interface will be used. (To make the examples shorter, the `Equals` method required to properly implement the interface will not be included.)

For example, for an entity named `Person` with properties `Name` of type `string` and `Age` of type `int`, the generated code should be similar to the one in Listing 1.

If an approach supports transforming C# code, not just generating it, the example will instead be transforming simple classes containing only fields with


```

1  using System;
2
3  class Person : IEquatable<Person>
4  {
5      public string Name { get; set; }
6      public int Age { get; set; }
7  }

```

Listing 1: Running example result

the right types and names into the form above. Such transformation is too simple to be useful in practice, but it is sufficient as a demonstration.

1.2 Manipulating C# source code

This section describes various approaches for generating and transforming C# source code. The resulting code then needs to be compiled by the C# compiler, as usual.

1.2.1 T4

Text Template Transformation Toolkit (T4) [7] is a tool for generating text by interspersing snippets of the text to generate with fragments of C# code to control how the text is generated. The resulting text can be in any language, including C#.

T4 does not have any special way of accessing other source code, which makes it most suitable for generating code based on external data. Its text-based nature gives it flexibility, but also makes using it fairly hard, since generating C# code effectively requires writing two interleaved programs, without any help from the IDE, because T4 integration into Visual Studio is very limited.

Example

The code to generate entities using T4, as required by the running example, can be seen in Listing 2.

The code highlights another issue with T4: indentation. It is hard to keep indentation of both the generated code and the generating code consistent, especially since any whitespace outside of T4 tags will be included in the output.

1.2.2 CodeDOM

Code Document Object Model (CodeDOM) [8] is a library for generating source code by using a language-independent object model. It is fairly easy to use, but is limited in what language features it supports, due to its language-independent nature and due to it not being updated since .NET Framework 2.0. Some of these limitations can be worked around by using string-based “snippet” objects, but using them means negating the advantages that CodeDOM has. Some examples of features it does not support are declaring `static` classes, LINQ query expressions or declaring auto-implemented properties.

```

1 <#@ assembly name="System.Runtime.dll" #>
2 <#@ assembly name="$(TargetDir)CSharpE.Samples.Core.dll" #>
3 <#@ import namespace="CSharpE.Samples.Core" #>
4 <#@ output extension=".cs" #>
5 using System;
6
7 <# foreach (var entityKind in EntityKinds.ToGenerate) { #>
8 class <# entityKind.Name #>
9     : IEquatable<<# entityKind.Name #>>
10 {
11 <# foreach (var property in entityKind.Properties) { #>
12     public <# property.Type #> <# property.Name #> { get; set; }
13 <# } #>
14 }
15 <# } #>

```

Listing 2: T4 example

Example

The code to generate entities for the running example using CodeDOM can be seen in Listing 3.

Especially notice that the code has to manually generate backing fields for properties, because CodeDOM does not support auto-implemented properties.

1.2.3 Roslyn

Roslyn [5] is the C# (and Visual Basic .NET) compiler, which can also be used as a library for manipulating C# code, including parsing, transformation and generation. Its object model was primarily designed to be used in the Visual Studio IDE, which is why it is very detailed (so it can accurately represent any source code, including erroneous or incomplete code) and also immutable (so that multiple IDE services can operate on the same model).

Roslyn contains several related Application programming interfaces (APIs) for manipulating source code, each useful in different situations:

- The **SyntaxTree** API forms the basis of Roslyn and represents only syntactic information about code. This means it can be used for just a single source file and makes it very efficient, especially when creating the **SyntaxTree** for code that contains only a small change relative to another **SyntaxTree**.

On the other hand, no semantic information is available from this API, so for example for the expression $F(A.B)$, it is not possible to determine whether F is a method or a delegate, whether A refers to a type or a variable, or whether B is a field, a property, or a method group.

The **SyntaxFactory** class can be used to create new nodes for this API.

- The **SemanticModel** class can be used to answer semantic questions about some part of a **SyntaxTree**.

The disadvantage of using this class is that it requires a full compilation, which includes all files in a project and also all of its dependencies. It is also less efficient, especially when a change is made.

```

1  var ns = new CodeNamespace();
2  ns.Imports.Add(new CodeNamespaceImport("System"));
3
4  foreach (var entityKind in EntityKinds.ToGenerate)
5  {
6      var entityType = new CodeTypeDeclaration(entityKind.Name);
7      entityType.BaseTypes.Add(new CodeTypeReference(
8          "IEquatable", new CodeTypeReference(entityKind.Name)));
9
10     foreach (var property in entityKind.Properties)
11     {
12         var propertyType = new CodeTypeReference(property.Type);
13
14         entityType.Members.Add(new CodeMemberField
15             {
16                 Name = property.LowercaseName, Type = propertyType
17             });
18
19         var fieldReference = new CodeFieldReferenceExpression(
20             new CodeThisReferenceExpression(), property.LowercaseName);
21
22         entityType.Members.Add(new CodeMemberProperty
23             {
24                 Attributes = Public | Final,
25                 Name = property.Name,
26                 Type = propertyType,
27                 GetStatements =
28                 {
29                     new CodeMethodReturnStatement(fieldReference)
30                 },
31                 SetStatements =
32                 {
33                     new CodeAssignStatement(fieldReference,
34                         new CodePropertySetValueReferenceExpression())
35                 }
36             });
37     }
38
39     ns.Types.Add(entityType);
40 }
41
42 var compileUnit = new CodeCompileUnit { Namespaces = { ns } };
43
44 using (var writer = new StreamWriter("Entities.cs"))
45 {
46     new CSharpCodeProvider().GenerateCodeFromCompileUnit(
47         compileUnit, writer, null);
48 }

```

Listing 3: CodeDOM example

The `SemanticModel` class surfaces semantic information in two forms:

- The `ISymbol` API can be used to get semantic information about members declared or referenced by a piece of code.

For example, for the expression `Console.WriteLine(42)`, this API could return the symbol for the `System.Console.WriteLine(int)` method. That symbol could then be used to find out more semantic information about that method, like the assembly it is contained in.

- The `IOperation` API is an alternative representation of statements and expressions as a language-independent abstract syntax tree. It includes semantic information in the form of `ISymbol` objects.

- The `SyntaxGenerator` class offers an alternative, language-independent way of generating Roslyn syntax nodes. It is part of the Workspaces layer of Roslyn, which means that using it on its own requires some additional setup. `SyntaxGenerator` has a more semantic view of code, which can be easier than generating the exact syntax using `SyntaxFactory`.

Since `SyntaxGenerator` is language-independent, it is using `SyntaxNode` in its API to represent any kind of syntax node, because `SyntaxNode` is the common base class for syntax node types in different languages. This approach makes `SyntaxGenerator` less type-safe when compared with `SyntaxFactory`, which uses specific `SyntaxNode`-derived types in its API.

Example

The code to transform entities for the running example using the `SyntaxTree` API can be seen in Listing 4.

Note that this code is heavily using `SyntaxFactory` to create new syntax nodes, but that type is not visible due to use of `using static`, to make the code more succinct.

Also notice how immutability makes transforming code harder by requiring the use of methods such as `ReplaceNodes` and how the level of detail leads to very verbose code, in some cases requiring even the specification of individual semicolons.

Code to transform entities for the running example using `SyntaxGenerator` can be seen in Listing 5.

Note that `SyntaxGenerator` does not support auto-implemented properties, so the example has to create properties with backing fields.

Also note that `SyntaxGenerator` generates overly verbose code (for example, `global::System.IEquatable<global::Person>`), but this is counterweighted by being able to use another Workspace-layer service, `Simplifier`, to make the code simpler.

TODO: More obscure libraries like RoslynDOM

```

1  var compilationUnit = ParseCompilationUnit(EntityKinds.ToGenerateFromSource);
2
3  compilationUnit = compilationUnit.ReplaceNodes(
4      compilationUnit.DescendantNodes().OfType<ClassDeclarationSyntax>(),
5      (_, classDeclaration) =>
6      {
7          classDeclaration = classDeclaration.AddBaseListTypes(
8              SimpleBaseType(QualifiedName(IdentifierName("System"),
9                  GenericName("IEquatable").AddTypeArgumentListArguments(
10                      IdentifierName(classDeclaration.Identifier)))));
11
12          var fields = classDeclaration.ChildNodes()
13              .OfType<FieldDeclarationSyntax>();
14
15          classDeclaration = classDeclaration.ReplaceNodes(fields,
16              (_, fieldDeclaration) =>
17              {
18                  var type = fieldDeclaration.Declaration.Type;
19                  var name = fieldDeclaration.Declaration.Variables.Single()
20                      .Identifier;
21
22                  return PropertyDeclaration(type, name)
23                      .AddModifiers(Token(PublicKeyword))
24                      .AddAccessorListAccessors(
25                          AccessorDeclaration(GetAccessorDeclaration)
26                              .WithSemicolonToken(Token(SemicolonToken)),
27                          AccessorDeclaration(SetAccessorDeclaration)
28                              .WithSemicolonToken(Token(SemicolonToken)));
29              });
30
31          return classDeclaration;
32      });
33
34  compilationUnit = compilationUnit.NormalizeWhitespace();
35
36  File.WriteAllText("Entities.cs", compilationUnit.ToString());

```

Listing 4: Roslyn SyntaxTree example

```

1  var project = new AdhocWorkspace().AddProject("MyProject", LanguageNames.CSharp)
2      .AddMetadataReference(
3          MetadataReference.CreateFromFile(typeof(object).Assembly.Location));
4  var document =
5      project.AddDocument("Entities.cs", EntityKinds.ToGenerateFromSource);
6  var g = SyntaxGenerator.GetGenerator(document);
7
8  var compilationUnit = (await document.GetSyntaxTreeAsync()).GetRoot();
9  var model = await document.GetSemanticModelAsync();
10 var compilation = model.Compilation;
11
12 compilationUnit = compilationUnit.ReplaceNodes(
13     compilationUnit.DescendantNodes().OfType<ClassDeclarationSyntax>(),
14     (_, classDeclaration) =>
15     {
16         SyntaxNode result = classDeclaration;
17
18         result = g.AddBaseType(result, g.TypeExpression(
19             compilation.GetTypeByMetadataName("System.IEquatable`1")
20             .Construct(model.GetDeclaredSymbol(classDeclaration))));
21
22         var fields = result.ChildNodes().OfType<FieldDeclarationSyntax>();
23
24         result = result.RemoveNodes(fields, default);
25
26         foreach (var fieldDeclaration in fields)
27         {
28             var type = fieldDeclaration.Declaration.Type;
29             var propertyName = g.GetName(fieldDeclaration);
30             var fieldName = propertyName.ToLowerInvariant();
31
32             var field = g.WithName(fieldDeclaration, fieldName);
33
34             var fieldAccess = g.IdentifierName(fieldName);
35             var property = g.PropertyDeclaration(
36                 propertyName, type, Accessibility.Public,
37                 getAccessorStatements: new[]
38                 {
39                     g.ReturnStatement(fieldAccess)
40                 },
41                 setAccessorStatements: new[]
42                 {
43                     g.AssignmentStatement(
44                         fieldAccess, g.IdentifierName("value"))
45                 });
46
47             result = g.AddMembers(result, field, property);
48         }
49
50         return result;
51     });
52
53 document = document.WithSyntaxRoot(compilationUnit.NormalizeWhitespace());
54 document = await Simplifier.ReduceAsync(document);
55
56 File.WriteAllText(
57     "Entities.cs", (await document.GetSyntaxRootAsync()).ToString());

```

Listing 5: Roslyn SyntaxGenerator example

1.3 Manipulating IL

This section describes various libraries that can be used for generating and transforming IL code.

Since IL is primarily meant to be produced by compilers and consumed by .NET runtimes, it is not a very convenient language for programmers.

1.3.1 System.Reflection.Emit

System.Reflection.Emit [9] is a library that can be used to generate an assembly at the IL level in memory, and then either directly execute code from that assembly or save the assembly to disk.

Example

Code to generate entities for the running example using Reflection.Emit can be seen in Listing 6.

Notice how generating properties requires understanding that they are represented using special methods and that generating bodies for those methods requires understanding the IL stack machine.

1.3.2 Mono.Cecil

Mono.Cecil [10] is a library that can be used for generating and transforming assemblies. It was written as part of the Mono project.

The main difference between Cecil and Reflection.Emit is that Cecil can be used to read assemblies, including their IL, and to modify them, while Reflection.Emit can only be used to create brand new assemblies. Another difference is that Cecil has its own type system, independent from the type system of the .NET runtime. This means that it can be used for example to work with assemblies that target a newer version of .NET than the currently executing one, or to work with assemblies that target an incompatible .NET implementation.

Cecil does not have an example of its usage shown here. While its APIs is different from Reflection.Emit, those differences are not relevant to this work.

1.4 Other approaches

While with the approaches mentioned in previous sections, then generated code (either C# or IL) fairly closely corresponds to the generating code, other options are possible. This section describes one such library.

1.4.1 Expression trees

Expression trees, [11] contained in the System.Linq.Expressions namespace, offer another representation of code.

Expression trees were first introduced in .NET Framework 3.5, to support translating of Language Integrated Query (LINQ) queries to existing query languages, such as Structured Query Language (SQL). This initial version included

```

1  var assemblyName = "MyAssembly";
2  var assembly = AssemblyBuilder.DefineDynamicAssembly(
3      new AssemblyName(assemblyName), AssemblyBuilderAccess.Save);
4  var module =
5      assembly.DefineDynamicModule(assemblyName, assemblyName + ".dll");
6
7  foreach (var entityKind in EntityKinds.ToGenerate)
8  {
9      var type = module.DefineType(entityKind.Name);
10
11      type.AddInterfaceImplementation(
12          typeof(IEquatable<>).MakeGenericType(type));
13
14      foreach (var propertyInfo in entityKind.Properties)
15      {
16          var propertyType = Type.GetType(propertyInfo.Type);
17
18          var field = type.DefineField(
19              propertyInfo.LowercaseName, propertyType, FieldAttributes.Private);
20
21          var property = type.DefineProperty(
22              propertyInfo.Name, PropertyAttributes.None, propertyType, new Type[0]);
23
24          var getMethod = type.DefineMethod("get_" + propertyInfo.Name,
25              MethodAttributes.Public | MethodAttributes.SpecialName,
26              propertyType, new Type[0]);
27
28          var il = getMethod.GetILGenerator();
29
30          il.Emit(OpCodes.Ldarg_0);
31          il.Emit(OpCodes.Ldfld, field);
32          il.Emit(OpCodes.Ret);
33
34          property.SetGetMethod(getMethod);
35
36          var setMethod = type.DefineMethod("set_" + propertyInfo.Name,
37              MethodAttributes.Public | MethodAttributes.SpecialName,
38              typeof(void), new[] { propertyType });
39
40          il = setMethod.GetILGenerator();
41
42          il.Emit(OpCodes.Ldarg_0);
43          il.Emit(OpCodes.Ldarg_1);
44          il.Emit(OpCodes.Stfld, field);
45          il.Emit(OpCodes.Ret);
46
47          property.SetSetMethod(setMethod);
48      }
49
50      type.CreateType();
51  }
52
53  assembly.Save(assemblyName + ".dll");

```

Listing 6: System.Reflection.Emit example

only expression-like constructs and the C# language supported compiling of expression lambdas to code that creates the corresponding expression tree.

In .NET Framework 4.0, expression trees were expanded with statement-like constructs, such as blocks, assignments or loops, to support the Dynamic Language Runtime (DLR). The result is a “language” that is still expression-based, which means that even constructs such as blocks or loops can have a result. This is somewhat similar to functional languages such as F#, which also do not differentiate between statements and expressions. The C# language was not updated to support the new constructs when translating lambdas to expression trees.

An expression tree can be inspected, often in order to be translated to some query language, or it can be executed. Depending on circumstances, executing expression trees is either done by using an interpreter, or by compiling them to IL using `Reflection.Emit` and then executing the result.

Expression trees can only represent expressions and statements, not types or their members, which limits their usefulness when generating code. This also means the running example is not applicable to expression trees.

1.5 Metaprogramming systems

TODO: More obscure tools, especially those that work with C#

1.5.1 PostSharp

PostSharp [12] is a commercial tool for transforming built assemblies at the IL level. It focuses on aspect-oriented programming (AOP), which is the idea that cross-cutting concerns (related pieces of code that are spread over the program, such as logging or code related to thread safety) should be specified separately from the rest of the code.

An aspect is applied to the target program element, such as a method, by attaching a specific attribute to it. The attribute can also be attached to a container, such as a type or an assembly, which applies the aspect to all relevant program elements in that container. This is called “attribute multicasting”.

PostSharp includes many built-in aspects and also allows specifying custom aspects. Custom aspects work by calling a user-defined method at a specific point, such as at the start of a method or before another method is called. The user-defined method can also be provided additional information about the modified method, such as its name or arguments.

This approach makes writing custom aspects easy, but it also means they are limited in what they can do and can have some performance penalty (due to allocation of the object that contains information about the modified method).

Like other tools that work at the IL level, PostSharp is also limited when it comes to changes to the shape of types, because any changes it makes will not be visible at compile time, only at runtime.

1.5.2 Fody

Fody [13] is an open-source tool for transforming the IL of assemblies. There are many published “add-ins” for Fody, usually distributed through NuGet, and custom add-ins can be written by modifying the assembly using the Mono.Cecil API (see Section 1.3.2). This makes custom add-ins hard to write, but it means they can perform any transformation. Though the limitations of IL-based tools still apply: changes to the shape of type will not be visible at compile time.

1.5.3 F# type providers

The F# language contains a feature called “type providers”, [14] meant for easier access to data sources. Type providers generate types at design-time (i.e. while editing code in an IDE) based on their input parameters and on usage of the generated types. These types can be either regular types that still exist after compilation (type providers using this approach are called “generative”) or their use can be transformed into some other code (“erased type providers”).

Type providers use “code quotations” to express code to execute. Code quotations serve a similar purpose in F# as expression trees do in C#.

2. Analysis

The goal of this work is to create a system for extending the C# language. It should be possible to use it create a wide variety of extensions, including:

- Extension similar to an F# type providers.
- Extension similar to a PostSharp aspects.
- Extension for entity types, which can generate constructors, members required for comparison and any other common boilerplate code.
- Extension that modifies how existing language feature operates, for example, improving the time complexity of recursive iterators from quadratic to linear.
- Extension that can be used to write a single method that performs numeric computation with any numeric type. This is easy to achieve with C++ templates, but impossible with C# generics, because they don't have a way of specifying operators required by the method.
- Extension that optimizes LINQ to Objects queries into efficient imperative code.
- Extension that converts an array of structures (AoS) into a structure of arrays (SoA), to improve performance of memory accesses in some cases.

TODO: More use cases from existing code generators?

Writing these extensions should be fairly easy and using extensions should not cause performance issues at design-time, build-time or run-time.

Such a system will require two major parts: an API for representing and modifying code used by extensions; and a component that drives extensions by applying their transformations at the appropriate time.

2.1 Representing code

The basic choices for representing C# code in an API are: as C# code, as IL, as some other form.

IL can be ruled out, because it is hard to use due to its low-level nature, especially when it comes to features like `async-await` (the C# compiler transforms `async` methods into state machines at the IL level).

Using a custom form would effectively require creating a new programming language (though one that does not necessarily have a textual form). The main disadvantage of doing that is that users would have to learn the new language and so it is generally the right choice only when no existing language is suitable.

This leaves the last option: using C#. This approach ensures that extension authors are already familiar with the used language, they only need to learn the API used to represent it. It also means that the output of the system will be in

C#, so existing tools for C# can be used. One disadvantage is that extensions can only use features available from C#. For example, the `calli` IL opcode is out of reach.

Putting this all together: C# is the best choice for forming the basis of the API for representing code for this system.

Now that we know that the API will represent C# code, we need to decide how exactly it should look like:

- The API should be a .NET library.

It would be possible to use an existing transformation language (such as Extensible Stylesheet Language Transformations (XSLT)) or create a new one. An existing language might not suit well the needs of transforming C# code and it would also require fitting C# code into the language's model (Extensible Markup Language (XML) in the case of XSLT). A new language would be unfamiliar to users and would lack tooling, at least initially. This means it would have to provide significant benefits to make creating it worth it.

The C# language has sufficient capabilities to express the transformations required to implement C# extensions and it is guaranteed to be familiar to the target group of this system: C# programmers. This means making the API a .NET library usable from C# is the best option.

- The API should be mutable.

Immutable persistent APIs (such as the one used by Roslyn) are useful when multiple versions of the same object need to be preserved (for example, for the “Undo” button in a code editor) or when multiple threads need access to the same object. Their disadvantage is that they make modifying objects harder: any change to a leaf of an object tree needs to be propagated to the root of the tree, creating new objects along the way.

This system does not need to keep multiple versions of objects, but it might be useful to parallelize it. For example, when two extensions modify different sections of code, it could be advantageous to execute them in parallel. Nevertheless, because of the focus on ease of use, a mutable API is the better choice.

- The API should not reflect the syntax of C# too closely.

In contrast with Roslyn, this API does not need to be able to represent, preserve and generate every syntactic nuance of C#, though it has to ensure that semantics of code is not changed. The basic examples of this are whitespace and comments.

A more advanced example is the difference between declaring variables together (`int i, j;`) or separately (`int i; int j;`). The API could represent both syntactic forms the same: as a sequence of two variable declarations.

Another example are parentheses in expressions. They are useful in the (infix) textual representation to change or emphasize the order of operations. But when representing an expression as a tree of objects, the order

of operations is clear from the structure of the tree, so parentheses are not necessary.

- The API should respect the syntactic structure of C#.

In contrast with the previous point, the API should not be completely divorced from the syntax of C#. For example, the general structure of a simple method declaration in C# is: modifiers, return type, method name, parameters, method body. If possible, the API should follow the same order.

- The API should make common code simple.

In the previous point, the list of elements of a method declaration was not complete: method declarations can also have attributes, type parameters and constraints. But many methods will not have these optional elements, and it should not be required to explicitly specify that a method does not have some of these elements.

Another example are method arguments. Method argument is commonly just an expression, but it can also have a modifier like `ref` or a name. But generating a method call without these optional elements is likely going to be the most common case, so it should be simple and not burdened by the requirements of more complex cases.

- The API should be succinct.

The structure of real code is often complex, so the API should handle generating such code. Because of this, it should avoid any unnecessary repetition, such as CodeDOM's `Code` prefix, Roslyn's `Expression` suffix or even repeated use of the `new` operator. Roslyn's `SyntaxFactory` with its `static` methods serves as a fairly good model here: when combined with `using static`, it leads to code that does not repeat itself much.

At the same time, the API should not be too succinct by abbreviating names, for example the way the C function `strpbrk` is named. This leads to names that are hard to understand and that also violate Framework Design Guidelines. [15]

- The API should seamlessly include semantic information.

Semantic information can be very useful, so it should be easily accessible. There shouldn't be a barrier similar to Roslyn's `SemanticModel`, where syntactic information is included in a syntax tree, but semantic information has to be accessed separately.

On the other hand, accessing semantic information is more expensive in terms of performance than information based purely on syntax, so it might make sense to somehow discourage their use.

- The API has to be capable of handling invalid code.

For some extensions, it will be useful if its users can write code that is not valid C#, which will then be transformed by the extension to make it valid. For example, an extension that automatically implements an interface could

require that its users specify that interface in the list of interfaces a class implements, but then omit any implementation. That is not valid C# code, but it will be filled out by the extension.

Another reason is that extensions need to work at design-time, while the code is being edited, so that auto-completion can include members produced by extensions. This is especially important for extensions similar to F# type providers, where generating new members is the reason why they exist.

Note that not all invalid code has to be handled equally well. Specifically, it is not clear how to parse or represent code that is not syntactically valid C#, such as code that attempts to use an operator that does not exist in C#. Such code still has to be handled by the API, but not necessarily consistently and it could include error nodes, or some similar representation for errors.

On the other hand, for code that is syntactically valid, but semantically invalid, such as the example of interface with no implementation mentioned above, it is clear how to parse such code and so its representation should be consistent and should not include any errors.

- The API should not be language-independent.

Several existing APIs for representing C# code are language-independent, at least to some degree. But the goal of this work is only to make the C# language extensible, so doing this is outside its scope.

Since none of the existing APIs satisfy these criteria well, it will be necessary to create a custom API for this system.

An API that follows the principles explained above will work the best for its designed purpose: writing extensions for this system. But it could be also used for others purposes, with some limitations. For example, if such API was used to write a Roslyn analyzer, it should be able to detect semantic issues (“Was a disposable object correctly disposed?”) but would likely have problems detecting syntactic or stylistic issues (“Does the code use `int` and not `Int32`?”).

2.2 The system

The next step is to consider how the overall system of executing extensions and applying their transformations should work.

- The system has to have a design-time component.

The primary purpose of some of the possible extensions mentioned at the start of this chapter is to generate code that is meant to be directly accessed by the extension’s user, often generated in response to other parts of the user’s code. This means that the generation has to be performed while the user is editing their code, in other words, at design-time.

- The system has to have a build-time component, which should be separate from the design-time component.

All of the mentioned possible extensions need to modify the build output, so a build-time component is clearly necessary.

And since the design-time component is likely going to be tied to a specific IDE or code editor, while the build-time component should work in a variety of situations, like building from the command line or on a build server, the two components should be separate.

- The system should support extensions with different design-time generation requirements.

Possible extensions have various requirements on code generated at design-time and at build-time, and the system should handle all of them. These include:

- An extension that generates different code at design-time and at build-time.

An example of such extension is one that is similar to an erased type provider: At design-time, it generates members with no implementation, which are then accessed by the user. At build-time, the members generated at design-time don't exist and user code that uses them is transformed into some other form.

This effectively requires writing two different transformations, one for each stage.

- An extension that generates the same code at design-time and at build-time.

An example is an extension similar to a generative type provider: Members are generated at design time and the same members are still used at build time.

There is still a difference between the two stages: it is not necessary to generate implementation of generated members at design-time, which is especially useful since design-time transformations are more time-sensitive. But it shouldn't have to be required to write two similar transformations for this.

- An extension that generates no code at design-time, only at build-time.

An example is an extension similar to an aspect: The extension is activated by attaching an attribute to a code element. The attribute does not change, so it does not have to be generated and can come from a regular library. This means that no code has to be generated at design-time. At build-time, the relevant code is then transformed based on what the aspect does.

- The system should regenerate only code that depends on modified code at design-time.

Performance of the design-time component is important, because it directly affects user experience when editing code that uses extensions, especially when it comes to auto-completion.

To help with that, we can take advantage of the fact that when the user is editing their code, they usually only change one piece of it at a time. And

since parts of an extension's transformation usually only depend on specific pieces of user code, it should be possible to execute only the parts of the transformation that depend on changed pieces of user code.

Another reason why this should be done is that extensions can affect performance of the whole system in unpredictable ways and limiting how much extension code runs also limits that unpredictability.

Doing this might require extending the API of the system.

- The system should make experimenting when creating extensions easy.

For example, if the system included the API suggested in the previous point, using it makes the extension more efficient but also more complicated. Because of that, use of this API should be recommended, but optional. This way, experimenting with writing extensions or creating personal extensions is simple thanks to the simple API, while production extensions can be efficient thanks to the complex API.

- The system should allow distributing extensions through NuGet.

NuGet is an established distribution channel for regular .NET libraries, and also for other kinds of libraries, like Roslyn analyzers or Fody add-ins. This makes NuGet a good fit for distributing extensions for this system.

- The system should allow extensions to report errors and warnings about code.

Many extensions are likely going to have ways of using their API (generated or not) in ways that are suspicious or outright incorrect. The system should let extensions report these issues to the user, including identifying the problematic part of their code.

As a side-effect, it would be possible to write an extension that does not perform any transformations, it only reports errors or warnings. Such extension would serve the same purpose as a Roslyn analyzer and would have similar limitations than those mentioned at the end of the previous Section.

3. Design

Before starting actual design of the system, its name should be decided. This name would be used in names of namespaces, assemblies, NuGet packages and so on, so it should fit well with their requirements and conventions. The name should also be reasonably unique, easy to remember and not too long. The name chosen based on these principles is “CSharpE”, meaning “C#, extensible”.

As explained in the previous chapter, this system has two main tasks: representing code and transforming code. This means it is natural to split the project into two main parts: `CSharpE.Syntax` and `CSharpE.Transform`, respectively.

3.1 CSharpE.Syntax

The `CSharpE.Syntax` namespace contains all the types necessary for representing and modifying C# code starting from the project level and going down all the way to the expression level. There is no representation for solutions, because extensions work at the project level, which means solutions are not necessary.

Listing 7 shows inheritance hierarchy of this namespace.

3.1.1 General principles

There are some rules that apply though all levels of this API:

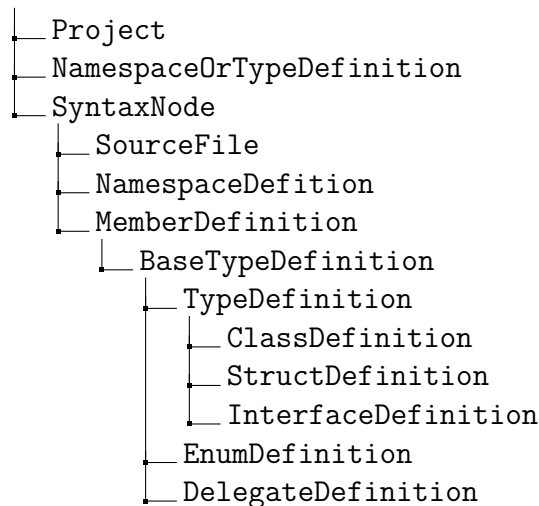
- Types in this API are regular mutable classes.
- Types that represent nodes in the C# syntax tree inherit from the common base class `SyntaxNode`. This includes most types from the source file level down.

A syntax node can have only one parent node, to make sure the syntax tree is actually a tree and so that mutating a node does not affect another, seemingly unrelated part of the syntax tree. But this parent node is not exposed in the API, for reasons explained in Section 3.2.

Syntax nodes can be deep cloned by calling the `Clone` method, which is a generic extension method. It is that way to avoid having to cast its result to the correct type, which would be necessary if `Clone` was a simple instance method on `SyntaxNode`. Syntax nodes are also cloned instead of assigning a new parent node, to maintain the tree shape.

- Collections of nodes are usually exposed as `ICollection<T>`. This interface is the most flexible out of the commonly used collection interfaces in .NET. Using an interface means that the user is shielded from the implementation detail of which specific collection type is used.
- The API includes implicit conversion operators when appropriate. Specifically, they can be used to convert from a definition to a reference (e.g. from

`TypeDefinition` to `IdentifierExpression` TODO: why not `NamedTypeReference`?)



Listing 7: Inheritance hierarchy of commonly used types in the CSharpE.Syntax namespace

or from a node to a simple wrapper for that node (e.g. from `Expression` to `ExpressionStatement`).

Using implicit conversions makes code shorter, but also harder to understand, because it is an operation that is not visible in the code. For this reason, all implicit conversion operators have an alternative form, usually a constructor of the target type.

- There is a `SyntaxFactory` type, which exists to make creating syntax nodes more succinct, when combined with `using static`. For example, it means that to use the `this` keyword, it is possible to write just `This()`, instead of `new ThisExpression()`.

3.1.2 Projects

At the top of this API is the `Project` class, which represents a collection of C# source files and library references. It also contains helper methods for accessing types from all files within a project, for extensions that do not care about which file contains which type.

The `Project` class is also point of interoperation between this API and Roslyn: `Project` can be constructed from a `CSharpCompilation` and it also exposes a `CSharpCompilation` as a property.

3.1.3 Source files

Each C# source file in the project is represented as an instance of the `SourceFile` class. A source file has name, text and a list of members, which are namespace and type definitions.

The object-oriented way to model this list would be to have a common base class or interface shared by `NamespaceDefinition` and `BaseTypeDefinition` (see next section) and use that for items in the list. But because namespace and

type definitions do not share many members and there are only two of them, it is likely users would write code specific to each of the two. For this reason, the list instead contains instances of a **struct** named **NamespaceOrTypeDefinition**, which is effectively a discriminated union of the two types. This is an approach common in functional programming.

Since this design is not very user-friendly, it is worth considering other options. One of them would be not having namespace definitions part of the syntax tree, but instead make namespace a property of type definition, which is how namespaces are represented in IL and Reflection. The main issue with this approach is that it would diverge too much from the structure of C# code, which could be confusing to users. For this reason, this option was not chosen.

Notice that **using** directives are not exposed in the API at all, because they are managed automatically by the API. This makes the API easier to use, at the cost of preventing users from choosing which syntax should be used, which is consistent with the principles outlined in the previous Chapter.

3.1.4 Types

The types that can be defined in C# are classes, structs, interfaces, enums and delegates. Classes, structs and interfaces are very similar in that all three can be generic, have a list of base types and have various members. Also, especially when it comes to classes and structs, it could make sense fairly often to manipulate them in the same way.

This means that a reasonable design would be to have the types for classes, structs and interfaces inherit from a common base class and this base class, along with the types for enums and delegates, should inherit from another base class. The problem with this design is naming: there is no established name that would apply to classes, structs and interfaces, but not to enums and delegates. For this reason, the names chosen for the two base classes are **TypeDefinition** and **BaseTypeDefinition**.

3.1.5 Members

Class, struct and interface definitions can contain various kinds of members, namely fields, methods, properties, events, indexers, operators, constructors, destructors (also known as finalizers) and nested type definitions. Not all kinds of members are valid for all three kinds of types, but since invalid members might be useful to some extensions, the API still allows them.

All kinds of members can have modifiers. While in source, member modifiers can be defined in different order (for example, both **public static** and **static public** are valid modifiers for a method), this order does not make a difference. For this reason, this API represents all modifiers using a single flags enum, **MemberModifiers**.

Because manipulating flags enums using bitwise operators can be cumbersome, the API also includes helper read-write properties for most valid modifiers for each kind of member. The exception to this are access modifiers. Because a member can have only one kind of declared accessibility (which includes all the access modifiers on their own and also the combinations **protected internal**

and `private protected`), access modifiers are surfaced as a read-write property named `Accessibility`, along with a read-only property for each kind of declared accessibility.

3.1.6 Statements

3.1.7 Expressions

3.1.8 References

3.2 CSharpE.Transform

TODO: Not just API!

4. Implementation

4.1 Syntax trees

4.2 Transformations

4.3 IDE integration

4.4 Example extensions

5. Comparison with existing tools

5.1 Reading and writing code

5.2 Transforming code

6. Future work

Conclusion

Bibliography

- [1] Phillip Carter. *Tour of .NET. Microsoft Docs .NET Guide*. 22 May 2017. URL: <https://docs.microsoft.com/en-us/dotnet/standard/tour> (visited on 18 June 2018).
- [2] Phillip Carter. *.NET architectural components. Microsoft Docs .NET Guide*. 23 Aug. 2017. URL: <https://docs.microsoft.com/en-us/dotnet/standard/components> (visited on 17 June 2018).
- [3] Ecma International. *Standard ECMA-335: Common Language Infrastructure (CLI)*. June 2012. URL: <https://www.ecma-international.org/publications/standards/Ecma-335.htm> (visited on 17 June 2018).
- [4] *C# 6.0 draft language specification. Microsoft Docs C# Guide*. 22 May 2018. URL: <https://docs.microsoft.com/en-us/dotnet/csharp/language-reference/language-specification/> (visited on 17 June 2018).
- [5] *The .NET Compiler Platform. GitHub*. URL: <https://github.com/dotnet/roslyn> (visited on 17 June 2018).
- [6] *NuGet Gallery*. URL: <https://www.nuget.org/> (visited on 30 July 2018).
- [7] Genevieve Warren. *Code Generation and T4 Text Templates. Microsoft Docs Visual Studio documentation*. 4 Nov. 2016. URL: <https://docs.microsoft.com/en-us/visualstudio/modeling/code-generation-and-t4-text-templates> (visited on 9 Aug. 2018).
- [8] Ron Petruscha. *Dynamic Source Code Generation and Compilation. Microsoft Docs .NET Framework Guide*. 30 Mar. 2017. URL: <https://docs.microsoft.com/en-us/dotnet/framework/reflection-and-codedom/dynamic-source-code-generation-and-compilation> (visited on 9 Aug. 2018).
- [9] Ron Petruscha. *Emitting Dynamic Methods and Assemblies. Microsoft Docs .NET Framework Guide*. 30 Aug. 2017. URL: <https://docs.microsoft.com/en-us/dotnet/framework/reflection-and-codedom/emitting-dynamic-methods-and-assemblies> (visited on 9 Aug. 2018).
- [10] *Mono.Cecil. Mono documentation*. URL: <https://www.mono-project.com/docs/tools+libraries/libraries/Mono.Cecil/> (visited on 9 Aug. 2018).
- [11] Bill Wagner. *Expression Trees. Microsoft Docs C# Guide*. 20 July 2015. URL: <https://docs.microsoft.com/en-us/dotnet/csharp/programming-guide/concepts/expression-trees/> (visited on 10 Aug. 2018).
- [12] *PostSharp*. URL: <https://www.postsharp.net/> (visited on 13 Aug. 2018).
- [13] *Fody. GitHub*. URL: <https://github.com/Fody/Fody> (visited on 14 Aug. 2018).
- [14] Phillip Carter. *Type Providers. Microsoft Docs F# Guide*. 2 Apr. 2018. URL: <https://docs.microsoft.com/en-us/dotnet/fsharp/tutorials/type-providers/> (visited on 14 Aug. 2018).

- [15] Krzysztof Cwalina and Brad Abrams. *Framework Design Guidelines. General Naming Conventions. Microsoft Docs .NET Guide*. 22 Oct. 2008. URL: <https://docs.microsoft.com/en-us/dotnet/standard/design-guidelines/general-naming-conventions> (visited on 17 Aug. 2018).

List of Figures

List of Tables

List of Abbreviations

AOP	aspect-oriented programming
AoS	array of structures
API	Application programming interface
CLI	Common Language Infrastructure
CodeDOM	Code Document Object Model
DLR	Dynamic Language Runtime
IDE	integrated development environment
IL	Intermediate Language
JIT	just-in-time
LINQ	Language Integrated Query
SoA	structure of arrays
SQL	Structured Query Language
T4	Text Template Transformation Toolkit
VB.NET	Visual Basic .NET
XML	Extensible Markup Language
XSLT	Extensible Stylesheet Language Transformations

A. Attachments

A.1 First Attachment