

# Applying Graph Partitioning-Based Seeding Strategies to Software modularisation

Ashley Mann<sup>1</sup> Stephen Swift<sup>1</sup> and Mahir Arzoky<sup>1</sup>

Brunel University London, UB8 3PH, UK

{Ashley.Mann, Stephen.Swift, Mahir.Arzoky @brunel.ac.uk}

**Abstract.** Software modularisation is a pivotal facet within software engineering, seeking to optimise the arrangement of software components based on their interrelationships. Despite extensive investigations in this domain, particularly concerning evolutionary computation, the research emphasis has transitioned towards solution design and convergence analysis rather than pioneering methodologies. The primary objective is to attain efficient solutions within a pragmatic timeframe. Recent research posits that initial positions in the search space wield minimal influence, given the prevalent trend of methods converging upon akin local optima. This paper delves into this phenomenon comprehensively, employing graph partitioning techniques on dependency graphs to generate initial clustering arrangement seeds. Our empirical discoveries challenge conventional insight, underscoring the pivotal role of seed selection in software modularisation to enhance overall outcomes.

**Keywords:** Software Engineering · Heuristic Search · Software Modularisation · Graph Partitioning

## 1 Introduction

### 1.1 General Background

As software systems grow, maintenance becomes challenging for incoming engineers unfamiliar with the original code, often leading to the need for significant overhauls or discontinuation of extensive legacy systems. To ensure sustainable management, creating modular subsystems is crucial. Instead of portraying them as clusters in the source code, a more practical approach is representing dependencies in graph form. Mancoridis et al. define the software modularisation problem as arising from the exponential complexity of interconnected software module relationships within evolving systems. This is often approached as a heuristic-search-based clustering problem to identify optimal representations by clustering subsystems based on the strength of their relationships [26].

The escalating complexity, often addressed through evolutionary computation, is evident in various software implementations, including both single-objective [4, 17] and multi-objective [19, 27] approaches. Pioneering methodolo-

gies aim to enhance the structure of software systems. Optimisation of sub-systems extends to diverse attributes, such as classes, methods, and variables. Methodological advancements now consider type-based dependence analysis [23], multi-pattern clustering [9], and effort estimation [32]. These efforts explore pre-processing and post-processing improvements alongside optimisation strategies.

The modularisation of software, mainly through heuristic search and evolutionary computation methodologies, extensively incorporates graph theory and data clustering. Academic works commonly use graph representations of software systems, employing data clustering for nodes and implementing algorithms to assess cluster quality [3, 20, 37]. Despite graph creation not inherently enhancing software engineers' understanding of architecture structure, language-independent graphs can focus on specific relationships or entire systems [11, 30, 35]. Clustering arrangements can be portrayed through various methods, such as a one-dimensional vector, a two-dimensional cluster-based structure, or a one-dimensional constrained representation known as a restricted growth function, which, despite its constraints, exhibits distinctive properties [8]. Clustering arrangement measurement typically addresses cohesion and coupling, striving for optimal cohesion within clusters and minimal coupling between clusters, fostering the creation of clearly defined groups [3].

## 1.2 Motivation

A recent study investigates diverse representations for clustering arrangements and different starting points, shedding light on the search space of software systems [1]. The study highlights the list-of-lists representation as the most robust, emphasising its significance in problem-solving. Notably, the paper suggests that the starting point choice is inconsequential, as various representations converge towards similar outcomes regarding final fitness, especially one and two-dimensional list-based ones.

This paper is motivated by exploring converging results based on starting points. Our primary objective is to determine whether alternative starting positions can replicate or potentially improve previous findings. If diverse starting positions tend to converge toward a similar region in the search space, we aim to uncover the reasons behind this convergence. Is there a basin of attraction leading to a potential global optimum solution, or do these methods unintentionally get stuck in closely adjacent local optima?

In recent years, a discernible research gap has emerged in clustering arrangement representation and software graph representation. Additionally, up to the present time, there is a notable absence of publications in the field of software modularisation specifically dedicated to addressing the concept of starting points. While we recognise that meta-heuristics, such as Iterated Local Search [22], can generate seeded starting points based on previous experimental iterations, our reference pertains to the primary initial search, distinct from subsequent iterations.

Building upon this motivation, we aim to explore innovative approaches for generating starting points that surpass the performance of previous experiments.

If our findings suggest the existence of a basin of attraction, our goal is to devise more efficient methods to reach this point faster than conventional approaches. However, even if the evidence points in a different direction, our overarching objective is to develop a more efficient method for navigating and exploring the search space.

This paper focuses on enhancing software system clustering by integrating graph partitioning techniques with seeded search methods applied to graph-based representations. Situated within Search-Based Software Engineering, our research particularly centres on software modularisation. To achieve our goal, we begin with a domain background, introduce innovative concepts, outline our experimental procedure, and present our results.

## 2 Related Work

### 2.1 Bunch and Munch

Exploring software modularisation can be achieved using tools such as Bunch [24] and Munch [5] [6]. Bunch, developed by Mancoridis et al., combines a Steepest Ascent Hill Climbing (SAHC) and Genetic Algorithms for improved clustering arrangements [24,25]. On the other hand, Arzoky et al., Munch employs Random Mutation Hill Climbing (RMHC) for enhanced performance and ease of implementation [5]. Both strategies use different fitness functions - Bunch utilises the *MQ* fitness function, while Munch employs *EVM* and *EVMD* [5,25,34]. Despite employing different measurement strategies, MQ, EVM and EVMD yield similar clustering results [18]. However, the exhaustive nature of Bunch may hinder performance when runtime is a critical consideration.

### 2.2 Starting Points and Search Space

In the context of a heuristic-search-based clustering problem, the quest for optimal solutions necessitates delving into the search space, which comprises all conceivable arrangements of a clustering configuration. This exploration entails generating an initial clustering arrangement known as the starting point. Subsequently, through mutation (searching), this arrangement is modified and compared to the graph representation of the software. The goal is to enhance the clustering of nodes that demonstrate robust relationships. Before embarking on a search, a crucial decision lies in determining the optimal starting point for seeking an improved clustering arrangement.

Several starting points are available when searching for local optima, which, in our context, represents the nearest approximation to the optimal clustering arrangement that maximises the cohesion of each cluster within the search space. We provide three illustrative examples: we can cluster all nodes individually for maximum coupling (Figure 1), together for maximum cohesion (Figure 2), or randomly (Figure 3).

Fig. 1: Independent

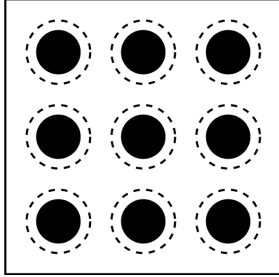


Fig. 2: All In One

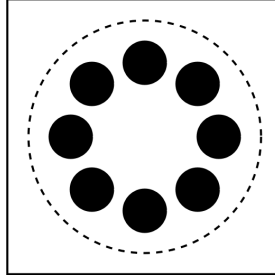
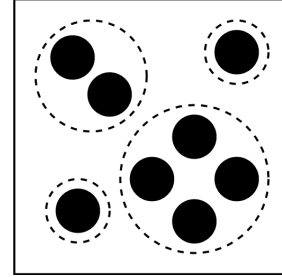


Fig. 3: Random



### 3 Research Questions

We aim to address research inquiries regarding our endeavour to discover improved starting points for software modularisation. We aim to uncover more effective strategies for achieving optimal outcomes. In this paper, we outline the following research questions that we intend to investigate:

1. What is the performance difference between graph-partitioned clustering arrangements and randomly generated ones when applied to large and small software systems?
  - (a) When using hill climbing with various initial clustering arrangements on the same software system, do the solutions converge to similar outcomes or do disparities persist?
  - (b) How do the runtimes of searches using graph partition and randomly generated clustering arrangements vary, and are there trade-offs between runtime and solution quality?
2. Is there a significant disparity between the Weighted Kappa<sup>1</sup> values of the final clustering arrangements and a gold standard<sup>2</sup>, and what is the nature of this comparison?

Initially, we aim to evaluate whether the graph-based initial clustering arrangements result in enhanced outcomes compared to randomly generated configurations through Munch. We aim to contrast the clustering patterns derived from graph partitioning with those generated randomly across software systems of varying sizes. Alongside assessing our fitness function, we analyse the documented improvements at the final iteration. This entails identifying the convergence point and scrutinising the runtime of the search, which encompasses both the initialisation of the starting configuration and the subsequent search process. By possessing information about the ultimate fitness value and the corresponding iteration when it is achieved, we aim to discern the genuine impact of the

<sup>1</sup> Weighted Kappa is employed to assess the similarity of clustering arrangements and is applied in sections 5.3, 6, and 7.

<sup>2</sup> A gold standard represents the theoretical best solution for a given problem, a rarity in real-world datasets where it is seldom known.

initial clustering configurations on the search dynamics. We aim to determine whether specific clustering arrangements contribute to a faster convergence, enabling us to refine our search methodology for reaching the convergence point earlier and mitigating the risk of potential time loss.

In addition to assessing the effectiveness of our initial clustering configurations based on fitness, convergence, and runtime, we also evaluate the final clustering arrangements against gold standards using Weighted Kappa (WK) [2]. WK serves as a measure of agreement between two clustering arrangements, explicitly focusing on modularisation. As the WK values increase, the level of agreement between the two solutions also rises. A WK value 1 signifies identical clustering arrangements, while 0 indicates empirical dissimilarity. A WK value of 0.5 or higher indicates a robust structural similarity between the two clustering configurations. We opt for WK over other methods, such as Adjusted Rand [29], due to its ease of implementation, longstanding presence in the field, and well-established interpretability/quality scale. The authors also note that WK and Adjusted Rand are identical.

## 4 Methods

Our focus now shifts towards the methodologies aligned with our exploration of optimal starting points for software modularisation. We present our selected search method and detail our implementation of graph partitioning designed to yield appropriate starting points.

### 4.1 Munch

As previously indicated, software modularisation can be characterised as a heuristic-search-based clustering problem. Therefore, our initial consideration lies in devising a strategy for heuristic search before delving into the discussion of our implementation of graph partitioning for generating starting points. We adopt a reverse-engineered adaptation of Arzoky et al.’s Munch to address this [6]. This adaptation has been enhanced to afford us the flexibility to determine the commencement of our exploration and the nature of our search strategy. We will now delve into an exploration of the various components that constitute Munch.

Foremost, Munch uses Module Dependency Graphs (*MDG*) as our graph-based representation of software systems. As defined by Mancoridis et al., *MDGs* illustrate subsystem connections to gauge relationships between components. In the context of our current research, we designate the nodes of *MDG* as software classes, and the edges represent interconnected relationships. *MDGs* prove versatile, capable of describing software structure over time or facilitating the segmentation of extensive software systems for enhanced comprehension. Let *MDG M* be an  $n$  by  $n$  symmetric binary matrix, where a 1 at row  $x$  and column  $y$  ( $M_{xy}$ ) indicates a relationship between software components  $x$  and  $y$ , and 0 indicates that there is no relationship. To avoid confusion throughout this paper,

MDG and graph are considered synonymous.

$$M_{xy} = \begin{cases} 1 & \text{if a relationship exists between } x \text{ and } y \\ 0 & \text{otherwise,} \end{cases}$$

For Munch, we adopt a list-of-list-based cluster representation based on its ease of implementation. A list-of-list clustering arrangement ( $C$ ) is defined as list  $([C_1, \dots, C_k])$ , with each subset list / cluster ( $C_i$ ) containing  $1, 2, \dots, n$  elements. These subsets must be non-empty ( $C_i \neq \emptyset$ ), and they should not share any common items ( $C_i \cap C_j = \emptyset$ ) for different subsets. Effective optimisation problem-solving requires consideration of the search space, exploration strategy, and fitness function. Equation 1 illustrates all possible ways to partition  $C_k$  clusters containing  $n$  elements. Note that  $1 \leq k \leq n$ . We justify opting for lists over sets in the implementation, emphasising the advantages of simpler implementation and reduced computational complexity, particularly in scenarios involving non-indexed sets. Since each cluster and cluster element requires indexing, the search space aligns with equation 1, deviating from the set nature characterised by Bell( $n$ ).

$$\sum_{k=1}^n \left( \frac{n!}{k! \cdot (n-k)!} \cdot k! \right) \quad (1)$$

Before delving into the search strategy of Munch, it is essential to define the fitness function. The primary goal of a search strategy is to uncover a clustering arrangement that most effectively aligns with the ideal modular structure of the software system. The assessment entails analysing the subsets  $[C_1, \dots, C_k]$ , where the elements ( $C_i$ ), representing  $1, 2, \dots, n$ , illustrate their relationships within the *MDG*. To avoid confusion, we will refer to the subsets as clusters.

For our replication of Munch, it is unsurprising we introduce *EVM* as our selected fitness function. We opt for *EVM* over Bunch's *MQ* due to its demonstrated robustness against noise and suitability for real-world software systems, as substantiated by research [18]. When provided with an arrangement  $C$  and an *MDG*, *EVM* evaluates and scores each cluster by considering the number of intra-relationships in the *MDG*. To prevent any potential confusion, we establish the definition of *EVM* as the aggregate of individual cluster scores, denoted as *SubEVM* (refer to Equations 2 and 3). *EVM* aims to maximise the score of relationships within a specified clustering arrangement. However, a potential drawback exists, as *EVM* may mistakenly assign high scores to clustering arrangements with high cohesion. Even minor adjustments to a solution can significantly enhance its fitness.

$$EVM(C, MDG) = \sum_{i=1}^k SubEVM(C_i, MDG) \quad (2)$$

$$SubEVM(C_i, MDG) = \sum_{a=1}^{|C_i|-1} \sum_{b=a+1}^{|C_i|} (2M(C_{ia}, C_{ib}) - 1) \quad (3)$$

To enhance the efficiency of Munch, we incorporate Arzoky et al.’s *EVMD*. This method generates a score aligning with *EVM* by integrating past *EVM* outcomes and determining the new result based on the classes designated for exchange. It demonstrates computational efficiency by computing the new fitness before implementing any modifications. Throughout this paper, we choose to utilise *EVM* as a collective term, encompassing both *EVM* and *EVMD*, to prevent potential confusion in future discussions about *EVM*.

Concerning the mentioned modifications, the inclusion of *EVMD* enables the execution of ‘Try/Do Swaps.’ This variant of Small-Change, involving the random mutation of clustering arrangements, reduces computational overhead by initially testing the result of a small change (Try Swap) before actual implementation (Do Swap). To effectively utilise *EVMD*, the swapping process is limited to two elements simultaneously.

Finally, our focus shifts to the heuristic search. As mentioned earlier, Arzoky et al.’s Munch primarily employ RMHC as its heuristic search method. Despite implementing the ability to alter the heuristic search in our Munch, we opt to persist with RMHC. This choice is motivated by its reliability, ease of implementation, and superior performance compared to stochastic heuristics, such as SAHC. Below, we present Algorithm 1, elucidating how Munch searches for enhanced clustering arrangements. For practical reasons, we choose to employ *EVM* in the pseudocode example, even though we leverage Arzoky et al.’s *EVMD* fitness function to enhance performance:

---

**Algorithm 1** Munch
 

---

```

1: procedure MUNCH(Iterations, MDG)
2:   Let C be a clustering arrangement           ▷ Random or Seed Starting Point
3:   Let F = EVM(C, MDG)                       ▷ Current Fitness
4:   for i = 1 to Iterations do
5:     Let C' = C                               ▷ Copy of C
6:     Choose two random clusters X and Y (X ≠ Y) from C', ▷ Move Operator
7:     Move a random variable from cluster X to Y in C'     ▷ Move Operator
8:     Let F' = EVM(C', MDG)                         ▷ New Fitness
9:     if F' ≥ F then                                   ▷ Compare Fitness
10:      Swap X and Y in C'                               ▷ Do Swap
11:      Let C = C', F = F'                             ▷ Continue using best Solution
12:     end if
13:   end for
14: end procedure                                       ▷ Output is C

```

---

## 4.2 Graph Partitioning

So far, we have established the importance of graphs and clustering arrangements regarding software modularisation. Now, we focus on using the structure graphs to discover new clustering arrangement starting points. Specifically, our focus shifts to the Fiedler Vector [13]. This vector is linked to the second smallest eigenvalue, the Fiedler Eigenvalue, of a Laplacian Matrix [14]. Denoted as

$L_{n \times n}$ , a Laplacian Matrix is defined as  $L = D - A$  where  $D$  represents the degree of matrix  $A$  which represents the connections between nodes [10]. In this context,  $A \equiv MDG$ . The Fiedler Vector is distinctive in its capability to enable a nearly perfect binary split of any given matrix. With this characteristic in mind, we have developed a tool that generates starting points through the recursive decomposition of graphs until no more Fiedler Vectors can be produced.

We generate a tree structure to facilitate the recursive decomposition of input graphs. The root of the tree is our input software graph and clustering arrangement. The clustering arrangement must begin with all nodes placed in a single cluster. This initial cluster will be subsequently split alongside the graph, ultimately leading to our final clustering arrangement, representing a fully decomposed software graph.

Leveraging our understanding of the Fiedler Vector, we identify the Fiedler eigenvalue of its attributed graph at each tree node, deduce its associated eigenvectors, and establish a well-balanced, binary-split graph partition. Simultaneously, we split the associated cluster with each partition, ensuring a one-to-one relationship between the subgraph's nodes and the associated cluster concerning the root MDG. This approach allows us to maintain traceability as we proceed with the decomposition. The new branches that emerge from the root node are reintroduced into a recursive function that continues to iterate until it identifies all possible partitions.

### 4.3 Starting Points

After generating a tree, we have two starting point approaches. Algorithm 2 illustrates the initial method for creating a "Leaf" arrangement. We gathered all leaf nodes from the tree, identified by their lack of children. Subsequently, we arrange these leaf nodes in ascending order based on SubEVM (see Equation 3) and then incorporate nodes with unique clusters into our clustering arrangement. We organise all leaf nodes in ascending order to prevent branches from becoming disconnected at different depths, possibly leading to duplicate values. In an ideal scenario, all leaf nodes, regardless of depth, should be unique, and therefore, we incorporate this logic for peace of mind.

Algorithm 3 exemplifies our alternative approach to constructing a clustering arrangement. In this method, we recursively traverse the tree, evaluating the cohesion of each node in comparison to its children. This ensures the creation of a clustering arrangement containing all unique values, emphasising the highest possible cohesion within the context of the MDG for a given tree. We refer to this starting point as our "Max" arrangement.

Apart from Leaf and Max, our modified version of Munch can generate clustering arrangements randomly distributed uniformly, denoted as "Random." Due to publication constraints, we abstain from delving into the intricacies of this method. In summary, a "uniformly distributed random" arrangement is defined by a clustering setup generated through the utilisation of Bell Numbers, Stirling Numbers of the Second Kind [21, 33, 36], and their interconnected relationships [12].

**Algorithm 2** BuildLeaf

---

```

1: function BUILDLEAF(root)
2:   Let leafNodes be a stack of all leaves of root
3:   Sort leafNodes in descending order of SubEVM           ▷ see Equation 2
4:   Let C be empty clustering arrangement
5:   while leafNodes is not empty do
6:     Let node be popped leafNodes                       ▷ Pop top node from stack
7:     if node is not in C then
8:       Let C = [C [node]]   ▷ Add unique cluster (node) to arrangement (C)
9:     end if
10:  end while
11:  return C
12: end function

```

---

**Algorithm 3** BuildMax

---

```

1: function BUILDMAX(root)                               ▷ Start Recursion
2:   Let C = empty clustering arrangement
3:   Let C = POPULATEARRANGEMENT(root, C)
4:   return C
5: end function
6: function POPULATEARRANGEMENT(node, C)               ▷ Recursive Function
7:   if node has children then
8:     Let Fp be SubEVM of parent node                   ▷ See Equation 2
9:     Let Fc be sum SubEVM of children                   ▷ See Equation 2
10:    if Fp > Fc then                                   ▷ If the SubEVM of parent node is greater
11:      Let C = [C [node]]                               ▷ Add node to arrangement C
12:    else                                                 ▷ Continue Recursion
13:      Let POPULATEARRANGEMENT(left child of node, C)
14:      Let POPULATEARRANGEMENT(right child of node, C)
15:    end if
16:  else
17:    Let C = [C [node]]                               ▷ Add leaf node to arrangement C
18:  end if
19: end function

```

---

## 5 Experimental Setup

Before presenting the Munch results of our graph partitioning tool, we need to establish an empirical framework.

### 5.1 Graph Collection and Pre-processing

First, we must collect software systems. Throughout our research, we developed a specialised tool that extracts open-source software systems using GitHub's RESTful API [15]. GitHub is our platform of choice for several compelling reasons. With a substantial user base exceeding 94 million developers, a continuously growing number of 52 million open-source repositories, and a cumulative

total of 413 million contributions [16], we have access to a wide and diverse range of graphs.

Collecting and forming these graphs is often neglected in academic literature, creating a challenge in determining the authenticity of these systems — whether they are genuinely open-source, artificially generated, or specific to certain industries. Generating MDGs requires understanding the relationships between each class within a given software system. This can be achieved using software metric tools such as SciTools Understand [31], which provide pairwise relationships to build a symmetric graph. After our extractor downloads the desired software system, we manually process each system using SciTools Understand. Future efforts will explore using GitHub’s TreeSitter parsing system [7] to automatically generate MDGs.

In this experiment series, we collect 50 "Small" open-source MDGs with class counts from 100 to 300, chosen based on relevance and high popularity ("stars") using the GitHub API. Due to storage constraints and the laborious manual MDG creation, we aim to develop an automated MDG generator, contemplating additional storage allocation pending study outcomes. Additionally, we have five "Big" MDGs (class counts: 1000 to 1500) sourced from prior research and industry collaboration, allowing exploration of size and characteristic-based result variations. Refer to Appendix A for a detailed breakdown. The terms 'Small 50' and 'Big 5' distinguish the two MDG groups in this paper.

## 5.2 Experiment Setup

Our experiments are described as follows. First, we collect Munch results for each MDG using all starting point combinations and iterations, as outlined below. Secondly, we collect Gold Standard results involving high-iteration/high-fitness outcomes to compare with our initial experiments. Finally, we analyse the results and present our findings concerning our outlined Research Questions.

1. For each experiment, for every graph (Small 50 and Big 5):
  - (a) Select one of three starting points (Leaf, Max, Random)
  - (b) Select one of three iterations (10k, 100k, 1m)
  - (c) Run Munch
  - (d) Document final iteration statistics and associated clustering arrangement
  - (e) Repeat Steps a-c 250 times
2. Repeat Step 1 until all starting points and iterations are explored.

For our gold standards, we generate a Random starting point for each graph and run Munch for 100 million iterations, collecting the same information as in our initial experiments. We repeat the process 250 times to ensure that we compare our initial experiment clustering arrangements to an absolute-best gold standard. Although conducting more iterations would have been preferable, it was impractical due to the extended runtime, taking several days per graph. To streamline experimental runs with our chosen iteration increments, we implemented parallel thread management, allowing multiple instances of Munch to run concurrently while optimising CPU and memory usage.

### 5.3 Data Collection and Analysis

We gather data on the fitness scores of the ultimate clustering configurations, pinpoint the convergence point (the last iteration demonstrating improved fitness), and gauge the runtime. Furthermore, we document the final clustering configurations into text files. Employing these files, we have crafted a bespoke tool to methodically evaluate the WK between the ultimate configurations derived from our initial points in contrast to our gold standards.

We have compiled a dataset of 275,000 files, combining the initial experiment results and gold standards. To enhance the manageability of these results for analysis, we employ MS Access, MS Excel, and Python for data processing. Due to the extensive volume of results and constraints in page space, our principal methodology involves computing averages across all data. Additionally, we streamline our findings by identifying and formatting the optimal results, providing a count of these instances per starting point type, thereby highlighting the suitability of each.

## 6 Results

Initially, we present RMHC results for each starting point category across selected iterations. Our research evaluates the performance of diverse starting points in searches across graphs of varying sizes, considering fitness, convergence, and runtime. The goal is to identify similarities or disparities in these aspects based on our predefined research questions. When implemented on large and small software systems, the performance differences among Leaf, Max, and Random are apparent in Tables 1 and 2<sup>3</sup>. Max consistently demonstrates superior fitness across iterations, as evidenced by the average final fitness values obtained from our three starting points over the specified iterations.

Table 3 details the average final convergence statistics across iterations, indicating the iteration where improvement was observed<sup>4</sup>. A strong resemblance between the average fitness and convergence strongly implies a correlation, potentially indicating a basin of attraction where all solutions converge. Compared to Random in Tables 1 and 2, Leaf and Max achieve final fitness levels more rapidly across all iterations, notably enhancing results. For smaller graphs, convergence is reached well before the considered iterations. Although final iterations align for smaller graphs, more iterations could enhance the likelihood of reaching local optima in larger datasets.

---

<sup>3</sup> Values highlighted in bold signify the highest average final fitness.

<sup>4</sup> Values highlighted in bold signify the shortest average convergence point.

Table 1: Average Final Fitness using Starting Point by Iterations

Size	10k			100k			1m		
	Leaf	Max	Random	Leaf	Max	Random	Leaf	Max	Random
102	60.66	<b>60.97</b>	60.52	<b>72.66</b>	73.03	72.54	<b>73.25</b>	73.24	72.98
105	<b>56.06</b>	56.02	56.05	<b>68.02</b>	67.70	67.84	<b>69.36</b>	69.19	69.34
112	88.61	<b>90.12</b>	87.53	99.51	<b>100.00</b>	99.92	101.04	<b>101.27</b>	101.11
119	66.96	<b>68.10</b>	66.20	80.10	<b>80.22</b>	80.17	<b>80.96</b>	81.05	80.90
123	99.52	<b>102.20</b>	98.07	<b>115.93</b>	115.69	114.90	117.81	<b>117.91</b>	117.66
127	82.34	<b>89.82</b>	79.45	105.68	<b>106.07</b>	105.27	107.22	<b>107.56</b>	107.04
129	81.80	<b>86.76</b>	80.50	98.54	<b>99.20</b>	98.85	99.80	99.86	<b>100.07</b>
130	37.61	<b>37.83</b>	35.59	53.02	<b>53.06</b>	53.20	53.83	53.66	<b>53.97</b>
135	97.75	<b>105.84</b>	96.86	120.68	<b>120.80</b>	120.64	122.20	122.05	<b>122.23</b>
135	101.48	<b>114.52</b>	100.76	133.83	<b>134.05</b>	133.91	<b>135.17</b>	134.92	135.00
141	74.43	<b>77.88</b>	71.39	88.22	88.55	<b>88.83</b>	89.83	<b>89.95</b>	89.93
145	<b>51.88</b>	51.86	47.85	<b>65.81</b>	65.62	65.67	<b>66.47</b>	66.38	66.37
151	105.95	<b>115.84</b>	102.02	<b>137.62</b>	137.01	137.22	139.28	139.13	<b>139.32</b>
158	84.08	<b>90.92</b>	76.42	112.66	112.53	<b>112.83</b>	115.26	115.29	<b>115.37</b>
161	85.58	<b>90.44</b>	82.20	<b>115.97</b>	114.97	115.75	118.91	118.85	<b>119.23</b>
163	83.98	<b>92.05</b>	73.40	110.25	<b>110.64</b>	109.73	<b>112.86</b>	112.78	112.19
164	79.77	<b>80.57</b>	75.16	105.46	<b>105.94</b>	105.41	107.88	107.93	<b>108.02</b>
169	42.46	<b>45.20</b>	35.43	73.16	<b>73.53</b>	73.21	75.32	<b>75.42</b>	75.37
172	86.92	<b>94.40</b>	82.73	123.58	122.60	<b>124.68</b>	<b>128.76</b>	128.23	129.08
172	100.64	<b>116.44</b>	91.55	137.20	<b>137.94</b>	137.17	<b>141.62</b>	141.37	140.93
172	15.80	<b>16.09</b>	11.80	<b>41.64</b>	41.49	41.37	44.28	44.39	<b>44.46</b>
175	80.77	<b>84.54</b>	71.78	<b>125.33</b>	125.20	125.11	133.86	133.88	<b>133.95</b>
175	108.70	<b>125.24</b>	101.23	164.62	<b>166.92</b>	164.51	171.25	<b>171.84</b>	171.03
176	82.36	<b>96.33</b>	72.36	112.78	112.84	<b>113.16</b>	115.76	<b>115.76</b>	115.69
179	101.70	<b>123.15</b>	93.20	148.36	148.16	<b>148.75</b>	151.90	151.34	<b>151.94</b>
198	98.01	<b>114.48</b>	84.57	139.20	<b>140.16</b>	139.19	144.09	144.28	<b>144.55</b>
199	90.33	<b>99.06</b>	83.39	153.60	150.50	<b>153.81</b>	166.76	<b>167.10</b>	166.66
200	91.70	<b>104.00</b>	70.86	128.12	<b>128.34</b>	127.91	131.67	131.27	<b>131.83</b>
206	113.88	<b>137.08</b>	103.77	182.22	<b>184.90</b>	181.95	191.01	<b>191.62</b>	190.96
208	62.53	<b>69.75</b>	52.00	127.38	<b>127.77</b>	127.51	130.66	<b>130.76</b>	130.72
209	82.90	<b>90.68</b>	59.67	131.20	<b>131.60</b>	129.78	137.77	<b>138.02</b>	137.24
210	65.22	<b>77.38</b>	48.58	102.34	<b>103.32</b>	101.98	106.34	<b>106.58</b>	106.40
210	76.32	<b>88.72</b>	54.56	116.98	<b>117.44</b>	116.69	118.22	<b>118.27</b>	118.06
211	66.62	<b>73.75</b>	44.26	113.07	113.32	<b>113.41</b>	120.87	120.88	<b>122.09</b>
214	118.44	<b>132.28</b>	100.06	187.69	<b>188.47</b>	187.41	198.19	<b>198.98</b>	198.08
216	43.36	<b>52.08</b>	47.34	120.72	<b>121.37</b>	120.73	126.76	126.76	<b>126.78</b>
224	110.88	<b>140.86</b>	89.71	184.60	<b>187.61</b>	184.31	197.92	<b>198.21</b>	198.14
233	92.68	<b>120.20</b>	75.74	173.61	<b>175.16</b>	172.78	<b>185.00</b>	184.81	184.60
234	89.56	<b>100.18</b>	57.73	143.93	<b>144.15</b>	142.66	<b>150.78</b>	150.76	150.65
234	111.31	<b>127.90</b>	84.24	184.22	<b>186.12</b>	183.92	<b>196.42</b>	196.26	195.99
235	93.87	<b>108.15</b>	65.81	158.05	<b>160.15</b>	158.59	172.26	<b>173.06</b>	172.60
240	82.87	<b>95.22</b>	54.92	148.99	<b>151.11</b>	147.78	159.10	159.13	<b>159.24</b>
240	126.10	<b>143.64</b>	98.22	<b>214.85</b>	213.21	213.49	234.84	<b>235.04</b>	234.47
242	108.40	<b>139.25</b>	76.39	177.71	<b>179.36</b>	177.09	186.90	<b>187.30</b>	187.28
246	72.88	<b>80.72</b>	41.40	139.55	<b>140.00</b>	139.34	148.82	148.52	<b>148.88</b>
252	99.48	<b>110.88</b>	61.12	159.87	<b>161.40</b>	159.00	174.66	<b>175.28</b>	174.62
252	134.83	<b>171.01</b>	108.42	243.30	<b>254.30</b>	242.45	279.19	<b>281.98</b>	278.47
258	123.15	<b>157.01</b>	94.90	208.43	<b>208.74</b>	206.76	<b>227.84</b>	225.93	227.14
264	114.19	<b>168.67</b>	72.40	227.13	<b>229.24</b>	226.55	239.49	238.79	<b>239.72</b>
294	131.71	<b>182.38</b>	75.39	251.48	<b>255.35</b>	249.92	276.91	276.69	<b>276.97</b>
Count	2	<b>48</b>	0	9	<b>34</b>	7	12	<b>20</b>	18

Table 2: Average Final Fitness using Starting Point by Iterations

Size	10k			100k			1m		
	Leaf	Max	Random	Leaf	Max	Random	Leaf	Max	Random
1037	170.92	<b>329.76</b>	-1360.34	345.94	<b>471.34</b>	123.48	731.89	<b>755.98</b>	737.26
1164	163.25	<b>276.05</b>	-1853.43	286.36	<b>376.87</b>	-188.73	588.53	<b>606.00</b>	589.68
1311	87.88	<b>225.91</b>	-2152.01	284.28	<b>400.46</b>	-297.58	681.85	<b>719.78</b>	675.82
1440	124.66	<b>283.65</b>	-2553.98	351.46	<b>479.48</b>	-520.69	741.54	<b>789.29</b>	708.47
1441	54.93	<b>146.11</b>	-2648.66	230.03	<b>304.93</b>	-718.30	543.32	<b>571.06</b>	493.34
Count	0	<b>5</b>	0	0	<b>5</b>	0	0	<b>5</b>	0

Table 3: Convergence Statistics

Iterations	Starting Point	Small 50			Big 5			Count
		Min	Max	Avg	Min	Max	Avg	
10k	Leaf	8970.44	9889.52	9560.28	9670.99	9887.30	9816.32	0
	Max	<b>8892.73</b>	<b>9873.90</b>	<b>9430.46</b>	<b>9634.59</b>	<b>9876.40</b>	<b>9799.48</b>	<b>6</b>
	Random	9129.95	9917.30	9710.81	9971.99	9977.35	9974.83	0
100k	Leaf	<b>53996.85</b>	96390.91	80512.02	<b>98831.28</b>	99120.06	98980.68	2
	Max	54158.65	<b>95036.25</b>	<b>79139.87</b>	98552.82	<b>99091.18</b>	<b>98786.13</b>	<b>4</b>
	Rand	54598.83	96333.32	80817.64	99769.87	99858.02	99828.65	0
1m	Leaf	84590.52	<b>691299.11</b>	375566.55	989895.56	993509.90	991965.96	1
	Max	<b>77164.98</b>	706030.09	<b>368672.25</b>	<b>987032.19</b>	<b>991801.38</b>	<b>989805.25</b>	<b>5</b>
	Random	85098.22	692030.05	377192.92	990564.15	994979.08	993397.87	0

In contrast to average final fitness and convergence, Table 4 highlights cumulative average runtimes presented for each start and subsequent search at various iterations, measured in milliseconds<sup>5</sup>. Notably, these reported runtimes represent summed average runtimes, excluding additional computational overhead related to data I/O. While Random clustering allows faster processing, the overall statistical significance of runtimes is debatable. This prompts consideration of potential trade-offs between runtime efficiency and solution quality.

Table 5 displays WK results, juxtaposing clustering configurations resulting from our initial starting points against gold standards<sup>6</sup>. In multiple statistics and iterations, Leaf and Max consistently surpass Random. The notable closeness between Max results and their corresponding gold standards in smaller graphs contrasts the generally low agreement observed for larger graphs.

Table 4: Average sum of runtime in milliseconds

Iterations	Small 50			Big 5		
	Leaf	Max	Random	Leaf	Max	Random
10k	298.797	332.881	<b>185.662</b>	24.231	<b>24.180</b>	199.044
100k	792.818	833.359	<b>408.770</b>	<b>119.465</b>	120.744	259.044
1m	5389.751	5607.426	<b>2574.516</b>	996.988	1012.166	<b>673.654</b>
Count	0	0	<b>3</b>	<b>1</b>	<b>1</b>	<b>1</b>

Table 5: WK against Gold Standard Statistics

Iterations	Starting Point	Small 50				Big 5				Count
		Min	Max	Avg	StDev	Min	Max	Avg	StDev	
10k	Leaf	0.215	0.565	0.386	<b>0.092</b>	0.002	0.002	0.002	<b>0.000</b>	2
	Max	<b>0.241</b>	<b>0.609</b>	<b>0.440</b>	0.094	<b>0.040</b>	<b>0.156</b>	<b>0.093</b>	0.043	<b>6</b>
	Random	0.198	0.580	0.358	0.102	0.010	0.027	0.015	0.007	0
100k	Leaf	0.398	<b>0.847</b>	0.637	<b>0.097</b>	0.012	0.021	0.016	<b>0.005</b>	3
	Max	0.376	0.835	<b>0.643</b>	0.099	<b>0.059</b>	<b>0.176</b>	<b>0.117</b>	0.044	<b>4</b>
	Random	<b>0.402</b>	0.843	0.635	0.097	0.058	0.107	0.077	0.022	1
1m	Leaf	0.490	0.897	0.715	0.084	0.128	0.317	0.198	0.079	0
	Max	0.480	<b>0.904</b>	<b>0.715</b>	0.086	<b>0.211</b>	0.359	0.264	<b>0.060</b>	<b>4</b>
	Random	<b>0.497</b>	0.897	0.715	<b>0.083</b>	0.197	<b>0.386</b>	<b>0.267</b>	0.084	<b>4</b>

<sup>5</sup> Values highlighted in bold signify the shortest runtime in milliseconds.

<sup>6</sup> Bold values signify the highest agreement.

## 7 Summary of Main Findings

In summary, we aimed to show that graph-partitioning can generate starting points capable of improving the results of software modularisation. We encapsulate the findings to address the research inquiries in the following manner:

- Max starting point:
  - Attains the highest average fitness over 10k, 100k, and 1m iterations, with a pronounced emphasis on lower iteration counts.
  - Attains the most average convergence across all iterations while sustaining the optimal average final fitness.
  - Attains the maximum average agreement (WK) with gold standards across 10k and 100k for both Small 50 and Big 5 graphs, highlighting noteworthy performance, especially in lower iterations.
- Leaf starting point:
  - Demonstrate fitness levels equal to or surpassing Random across all iterations, especially in the early stages
  - Surpasses Random with higher average fitness levels on large datasets at 10k and 100k iterations
  - Consistently exhibits faster convergence compared to Random.
- Random starting point:
  - Shows a quicker average total runtime in milliseconds compared to Leaf and Max.
  - Better suited for smaller datasets; however, an improvement over Max and Leaf necessitates higher iterations.
  - Demonstrates greater resemblance to the gold standard than Max in larger systems at 1m iterations.

Distinct fitness variations emerge among Leaf, Max, and Random, with Max consistently outperforming over 10k, 100k, and 1m iterations, notably in Small 50 vs. Big 5 comparisons. Random outperforms Max and Leaf at 1 million iterations for large datasets. However, Max proves to be more suitable for average fitness and faster convergence across iterations and graph sizes. Since there are currently no guidelines for determining the number of iterations based on the size or properties of an MDG, the most prudent approach would be to initiate seeding with Max before executing Munch. WK comparisons show Max starting points yield higher average agreements, with potential improvements around 70% and significant opportunities at 90% agreement in 1m iterations. Thorough exploration is vital for understanding software system graph intricacies. Our commitment to accelerating software modularisation drives deeper exploration, with partition-based clustering performing significantly, especially at smaller iterations, making it compelling for future software optimisation.

## 8 Generalisability

This publication focuses on utilising graph partitioning for software modularisation. However, the application of graph partitioning for optimising initial positions can extend to other graph-based optimisation problems, contingent on the chosen fitness function. Although we prioritize EVM for its simplicity, other alternatives like MQ are viable. Our aim is to inspire exploration of graph partitioning for seeded optimisation.

## 9 Future Work

We plan to integrate our graph-based initial clustering with metaheuristics, specifically incorporating seeded starting points into the history of Iterated Local Search, as part of our ongoing investigation [28]. This initiative seeks to evaluate the potential improvement in the exploration of the search space and overall efficiency. Furthermore, our goals include delving deeper into software systems' search space, exploring graph structure, convergence prediction, and other avenues for enhancing software modularisation.

## References

1. Intelligent systems and applications: Proceedings of the 2024 intelligent systems conference (intellisys) volume 1. In: *Lecture Notes in Networks and Systems #822*. p. 470. Springer (January 5 2024), author/s intentionally omitted for double-blind review
2. Altman, D.: *Skewed distributions. Practical statistics for medical research*. London, Chapman & Hall pp. 60–63 (1997)
3. Arasteh, B.: Clustered design-model generation from a program source code using chaos-based metaheuristic algorithms. *Neural Computing and Applications* **35**(4), 3283–3305 (2023)
4. Arasteh, B., Seyyedabbasi, A., Rasheed, J., M. Abu-Mahfouz, A.: Program source-code re-modularization using a discretized and modified sand cat swarm optimization algorithm. *Symmetry* **15**(2), 401 (2023)
5. Arzoky, M., Swift, S., Tucker, A., Cain, J.: Munch: An efficient modularisation strategy to assess the degree of refactoring on sequential source code checkings. In: *2011 IEEE Fourth International Conference on Software Testing, Verification and Validation Workshops*. pp. 422–429. IEEE (2011)
6. Arzoky, M., Swift, S., Tucker, A., Cain, J.: A seeded search for the modularisation of sequential software versions. *J. Object Technol.* **11**(2), 6–1 (2012)
7. Brunsfeld, M.: *Tree-sitter*, <https://github.com/tree-sitter/tree-sitter>, Accessed on 2023-11-01
8. Campbell, L.R., Dahlberg, S., Dorward, R., Gerhard, J., Grubb, T., Purcell, C., Sagan, B.E.: Restricted growth function patterns and statistics. *Advances in Applied Mathematics* **100**, 1–42 (2018)
9. Chen, Y.T., Huang, C.Y., Yang, T.H.: Using multi-pattern clustering methods to improve software maintenance quality. *IET Software* **17**(1), 1–22 (2023)
10. Chung, F.R.: *Spectral graph theory*, vol. 92. American Mathematical Soc. (1997)

11. Corradini, A., König, B., Nolte, D.: Specifying graph languages with type graphs. *Journal of Logical and Algebraic Methods in Programming* **104**, 176–200 (2019). <https://doi.org/https://doi.org/10.1016/j.jlamp.2019.01.005>, <https://www.sciencedirect.com/science/article/pii/S235222081730233X>
12. Devroye, L.: Sample-based non-uniform random variate generation. In: *Proceedings of the 18th conference on Winter simulation*. pp. 260–265 (1986)
13. Fiedler, M.: A property of eigenvectors of nonnegative symmetric matrices and its application to graph theory. *Czechoslovak mathematical journal* **25**(4), 619–633 (1975)
14. Fiedler, M.: Laplacian of graphs and algebraic connectivity. *Banach Center Publications* **1**(25), 57–70 (1989)
15. GitHub: Github advanced search (2023), <https://github.com/search/advanced>, Last Accessed on 23-11-01
16. GitHub: Octoverse 2022: 10 years of tracking open source (2023), <https://github.blog/2022-11-17-octoverse-2022-10-years-of-tracking-open-source/>, Last Accessed on 23-11-01
17. Gupta, N., Kumar, S., Gupta, V., Vijh, S.: Novel automatic approach using modified differential evaluation to software module clustering problem. *SN Computer Science* **4**(6), 816 (2023)
18. Harman, M., Swift, S., Mahdavi, K., Beyer, H.: An empirical study of the robustness of two module clustering fitness functions, pp. 1029 – 1036. *ASSOC COMPUTING MACHINERY* (2005), genetic and Evolutionary Computation Conference ; Conference date: 25-06-2005 Through 29-06-2005
19. Kang, Y., Xie, W., Wang, X., Wang, H., Wang, X., Li, J.: Mopisde: A collaborative multi-objective information-sharing de algorithm for software clustering. *Expert Systems with Applications* p. 120207 (2023)
20. Khan, M.Z., Naseem, R., Anwar, A., Haq, I.U., Alturki, A., Ullah, S.S., Al-Hadhrami, S.A., et al.: A novel approach to automate complex software modularization using a fact extraction system. *Journal of Mathematics* **2022** (2022)
21. L.H., H.: Stirling behaviour is asymptotically normal. *The Annals of Mathematical Statistics* **3**(2), 410–414 (1967)
22. Lourenço, H.R., Martin, O.C., Stützle, T.: Iterated local search. In: *Handbook of metaheuristics*, pp. 320–353. Springer (2003)
23. Lu, K.: Practical program modularization with type-based dependence analysis. In: *2023 IEEE Symposium on Security and Privacy (SP)*. pp. 1256–1270. IEEE (2023)
24. Mancoridis, S., Mitchell, B.S., Chen, Y., Gansner, E.R.: Bunch: A clustering tool for the recovery and maintenance of software system structures. In: *Proceedings IEEE International Conference on Software Maintenance-1999 (ICSM'99). 'Software Maintenance for Business Change'*(Cat. No. 99CB36360). pp. 50–59. IEEE (1999)
25. Mancoridis, S., Mitchell, B.S., Rorres, C., Chen, Y., Gansner, E.R.: Using automatic clustering to produce high-level system organizations of source code. In: *Proceedings. 6th International Workshop on Program Comprehension. IWPC'98* (Cat. No. 98TB100242). pp. 45–52. IEEE (1998)
26. Mitchell, B.S., Mancoridis, S.: Clustering module dependency graphs of software systems using the bunch tool. *Nat. Sci. Found., Alexandria, VA, USA, Tech. Rep* (1998)
27. Prajapati, A., Parashar, A., Rathee, A.: Multi-dimensional information-driven many-objective software remodularization approach. *Frontiers of Computer Science* **17**(3), 173209 (2023)

28. Ramalhinho-Lourenço, H., Martin, O.C., Stütze, T.: Iterated local search (2000)
29. Rand, W.M.: Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical association* **66**(336), 846–850 (1971)
30. Savić, M., Rakić, G., Budimac, Z., Ivanović, M.: A language-independent approach to the extraction of dependencies between source code entities. *Information and Software Technology* **56**(10), 1268–1288 (2014)
31. SciTools: Understand: The software developer’s multi-tool (2023), <https://scitools.com/>, Accessed on 2023-11-10
32. Tan, A.J.J., Chong, C.Y., Aleti, A.: Closing the loop for software remodularisation–rearrange: An effort estimation approach for software clustering-based remodularisation. *arXiv preprint arXiv:2303.06283* (2023)
33. Temme, N.M.: Asymptotic estimates of stirling numbers. *Studies in Applied Mathematics* **89**(3), 233–243 (1993)
34. Tucker, A., Swift, S., Liu, X.: Variable grouping in multivariate time series via correlation. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* **31**(2), 235–245 (2001)
35. Weiss, K., Banse, C.: A language-independent analysis platform for source code. *arXiv preprint arXiv:2203.08424* (2022)
36. Weisstein, E.W.: Stirling number of the second kind. <https://mathworld.wolfram.com/> (2002)
37. Yang, K., Wang, J., Fang, Z., Wu, P., Song, Z.: Enhancing software modularization via semantic outliers filtration and label propagation. *Information and Software Technology* **145**, 106818 (2022)

## 10 Appendix A

This Appendix showcases details about the software system MDG used in our experiments. Below, we showcase the following statistics for each software system:

1. ID
  - Each software system is assigned a unique identifier. We choose not to use the actual names of our software systems because our collection is sourced randomly from GitHub. These software system names can exhibit variation, and we intend to maintain professionalism and steer clear of potentially inappropriate names and software tools.
2. Nodes
  - Also known as vertices, these signify the number of software components (classes) within our Module Dependency Graphs (MDGs).
3. Edges
  - Denotes the number of relationships between software components.
4. Clustering Coefficient:
  - The extent to which nodes tend to cluster. A high score indicates a strong cohesion, while a low score indicates a higher coupling level. We present this statistic as these software systems exhibit remarkably low coefficients, indicating a high coupling level and a deficiency in the initial modular structure. There is potential here to investigate the nature of software structure over time, especially concerning the analysis of open-source software systems.

Table 6: 'Big 5' Software MDG Statistics

Identification	Nodes	Edges	Avg Degree	Clustering Coefficient
1	1037	5470	5.274	0.000
2	1164	2072	1.780	0.000
3	1311	5630	4.294	0.000
4	1440	4889	3.395	0.000
5	1441	3058	2.122	0.000

Table 7: 'Small 50' Open-Source Software MDG Statistics

Identification	Nodes	Edges	Avg Degree	Clustering Coefficient
01	102	312	0.061	0.007
02	105	257	0.047	0.002
03	112	436	0.070	0.007
04	119	343	0.049	0.002
05	123	440	0.059	0.004
06	127	409	0.051	0.004
07	129	357	0.043	0.002
08	130	225	0.027	0.001
09	135	387	0.043	0.002
10	135	510	0.056	0.004
11	141	333	0.034	0.001
12	145	274	0.026	0.001
13	151	595	0.053	0.002
14	158	422	0.034	0.001
15	161	413	0.032	0.001
16	163	414	0.031	0.001
17	164	686	0.051	0.002
18	169	387	0.027	0.001
19	172	609	0.041	0.002
20	172	591	0.040	0.002
21	172	357	0.024	0.000
22	175	737	0.048	0.004
23	175	749	0.049	0.003
24	176	371	0.024	0.000
25	179	467	0.029	0.001
26	198	552	0.028	0.001
27	199	1002	0.051	0.003
28	200	450	0.023	0.000
29	206	964	0.046	0.003
30	208	643	0.030	0.001
31	209	732	0.034	0.001
32	210	518	0.024	0.000
33	210	323	0.015	0.000
34	211	599	0.027	0.001
35	214	834	0.037	0.002
36	216	740	0.032	0.001
37	224	937	0.038	0.002
38	233	818	0.030	0.001
39	234	521	0.019	0.000
40	234	930	0.034	0.002
41	235	823	0.030	0.001
42	240	898	0.031	0.001
43	240	1115	0.039	0.002
44	242	836	0.029	0.001
45	246	672	0.022	0.001
46	252	1033	0.033	0.001
47	252	1591	0.050	0.004
48	258	1477	0.045	0.002
49	264	730	0.021	0.001
50	294	1275	0.030	0.001