# gpGeocoder Documentation

Documentation Overview                                                    release notes gpGeocoder-1.8.25.0-SNAPSHOT

## Result Classification

## Requirements

> Return a quality description of the results compared to the input.
> Describe classes for the result qualities.
> Describe characteristical features of the result.
> Allow a sorting of the results.
> Give Criteria for batch geocoding.

Explicitly **not** in the focus of the classification are the following issues:

> Give a continuous number to rate the result quality.
> Provide information about internals of the geocoder algorithm.

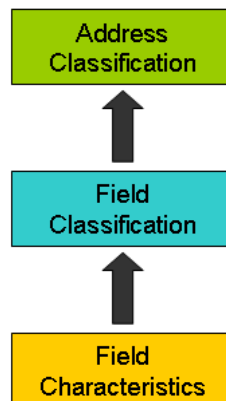From these requirements, the classification was specified as follows.

## Specification

### Introduction

The classification of a result address is a three-stage process:

> 1. **Field Characteristics** For each field all applicable characteristics are determined. The characteristics describe properties of a field, in the majority of cases applying comparison criterions on input and result.
> 2. **Field Classification** Based on only these characteristics, the classification of each single field is computed.
> 3. **Address Classification** The classification of the complete result address is derived from the classification of all fields.



A top-down description of each stage follows below:

### Address Classification

The classification of the complete address is done according to the rules that are given in the following table.

| Classification | Condition |
|---|---|
| **Exact** | (P == Exact \|\| NoInput) && <br> (C == Exact \|\| NoInput) && <br> (S == Exact \|\| NoInput) && <br> (H == Exact \|\| NoInput) |
| **Partially exact** | (P >= Partially exact \|\| NoInput \|\| NoResult) && <br> (C >= Partially exact \|\| NoInput) && <br> (S >= Partially exact \|\| NoInput) && <br> (H >= Partially exact \|\| NoInput \|\| NoResult) |
| **High** | (P >= High \|\| NoInput \|\| NoResult) && <br> (C >= High \|\| NoInput) && <br> (S >= High \|\| NoInput) && <br> (H >= Partially exact \|\| NoInput \|\| NoResult) |
| **Medium** | (((P >= Partially exact \|\| C >= Partially exact) && (S >= Partially exact \|\| NoInput)) \|\| <br> (P >= Partially exact \|\| NoInput \|\| NoResult) && (C >= Partially exact \|\| NoInput) && (S >= Partially exact \|\| NoInput \|\| NoResult) \|\| <br> (P >= Medium \|\| NoInput \|\| NoResult) && (C >= Medium \|\| NoInput) && (S >= Medium \|\| NoInput)) && <br> (H >= Partially exact \|\| NoInput \|\| NoResult) |

| Low | Rest |
|-----|------|

| Legend | |
|---|---|
| **P** | Postalcode |
| **C** | City/City2 |
| **S** | Street |
| **H** | House Number |

## Field Classification

| Classification | Availability | | | | Definition |
|---|---|---|---|---|---|
| | Postalcode | City/City2 | Street | HNr | |
| **Exact** | ✓ | ✓ | ✓ | ✓ | For details for each single field see below. |
| **Partially exact** | ✓ | ✓ | ✓ | ✓ | For details for each single field see below. |
| **High** | ✓ | ✓ | ✓ | ⊖ | Not contained in superior classes and Match quality high is set. |
| **Medium** | ✓ | ✓ | ✓ | ⊖ | Not contained in superior classes and Match quality medium is set. |
| **Low** | ✓ | ✓ | ✓ | ✓ | Not contained in any other class |
| **NoInput** | ✓ | ✓ | ✓ | ✓ | The input field is empty. |
| **NoResult** | ✓ | ✓ | ✓ | ✓ | The result field is empty and the input field is not empty. |

| Postalcode | |
|---|---|
| **Exact** | Match exact |
| **Partially exact** | Input is prefix of result |

| City/City2 + Street | |
|---|---|
| **Exact** | Match exact && Input is prefix of result |
| **Partially exact** | (Input is prefix of Multiword result \|\| Match except abbreviations \|\| Match except separators) && Match quality high && Result words include input |

| House Number | |
|---|---|
| **Exact** | Match exact |
| **Partially exact** | Match exact except additions |

## Field Characteristics

### Definitions

> The **definition** of a term is written bold letters. The *usage* of a defined term is written in italic letters.
> The steps listed in the column Preprocessing are applied in the given order on both input and result. Then, the characteristic is true if the preprocessed data fulfills the definition. Following preprocessing steps are available:
>> **SN** Special character normalisation: Special characters and umlauts are replaced (e.g ü->ue, é -> e, ...).
>> **CN** Case normalisation: Upper and lower case characters are treated in the same way.
>> **AN** Abbreviation normalisation: The replacements defined in the abbreviation dictionary are applied.
>> **XN** Affix normalisation: All affixes defined in the affix dictionary are removed from the string.
>> **PN** Phonetic normalisation: The phonetic replacements are applied.
>> **SepN** Separator normalisation: All separators (Space, -, /, ., etc.) are treated in the same way. Multiple consecutive separators are projected to a single one.
> Prefix: x is Prefix of another word y if x is not longer than y and characters on position i (i = 1 ... length( x )) in x and y are identical.
> Multiword Prefix: A list of words w is a Multiword Prefix of a list of words z, if an injective mapping from the words $w_1 ... w_n$ to the words $z_1 ... z_m$ exists such that each word $w_1 ... w_n$ is Prefix of its corresponding word in $z_1 ... z_m$. (Note that m >= n)

### Characterization bitfield values

The characterization is returned as a number that represents a bitfield. The values are given in the enumeration **eResultCharacteristics**.

| Characteristic | Availability | | | | Preprocessing | Definition | Examples | |
|---|---|---|---|---|---|---|---|---|
| | Postalcode | City/City2 | Street | HNr | | | Positive | Negative |
| **Input is prefix of result**<br><br>**rcInputIsPrefixOfResult** | ✓ | ✓ | ✓ | ⊖ | CN, SN, SepN | Input is *Prefix* of result | I: karl<br>O: Karlsruhe | I: karle<br>O: Karlsruhe |

| Classification | 1 | 2 | 3 | 4 | Codes | Description | Example A | Example B |
|---|---|---|---|---|---|---|---|---|
| **Input is phonetic prefix of result** <br> **rcInputIsPhoneticPrefixOfResult** | ✗ | ✓ | ✓ | ✗ | CN, SN, SepN, PN | Input is *Prefix* of result | I: karl <br> O: Karlsruhe <br> I: carl <br> O: Karlsruhe | I: karle <br> O: Karlsruhe |
| **Input is prefix of Multiword result** <br> **rcInputIsPrefixOfMultiwordResult** | ✗ | ✓ | ✓ | ✗ | CN, SN, SepN | Input is *Multiword Prefix* of result | I:Wei Beu <br> O:Beutelsbach Weihersdorf <br> I: Beu <br> O:Beutelsbach Weihersdorf | I:BeuWei <br> O: Beutelsbach Weihersdorf <br> I: Beu Wei <br> O:Beutelsbach |
| **Input is prefix of phonetic Multiword result** <br> **rcInputIsPrefixOfPhoneticMultiwordResult** | ✗ | ✓ | ✓ | ✗ | CN, SN, SepN, PN | Input is *Multiword Prefix* of result | I:Karls Turla <br> O:Karlruhe Durlach <br> I:Carls <br> O:Karlruhe Durlach | I:BeuWei <br> O: Beutelsbach Weihersdorf <br> I: Beu Wei <br> O:Beutelsbach |
| **Match except abbreviations** <br> **rcMatchExceptAbbreviations** | ✗ | ✓ | ✓ | ✗ | CN, SN, SepN, AN | Input is *Multiword Prefix* of result | I:Stumpfstrasse <br> O:Meine Stumpfstr | I:Stumpfstrasse <br> O:Stumpfweg <br> I:Stumpfstrasse <br> O:Stampfstr |
| **Match except affixes** <br> **rcMatchExceptAffixes** | ✗ | ✓ | ✓ | ✗ | CN, SN, SepN, XN | Input is *Multiword Prefix* of result | I:Stumpfweg <br> O:Stumpfstr | I: Stumpfstra <br> O: Stumpfstr |
| **Result words include input** <br> **rcResultWordsIncludeInput** | ✗ | ✓ | ✓ | ✗ | CN, SN, SepN, AN | Input is *Multiword Prefix* of result | I: Beu <br> O:Beutelsbach Weihersdorf | I: Beu Wei <br> O:Beutelsbach |
| **Input words include result** <br> **rcInputWordsIncludeResult** | ✗ | ✓ | ✓ | ✗ | CN, SN, SepN, AN | Result is *Multiword Prefix* of input | I: Beutelsbach Wei <br> O:Beutelsbach | I: Beu <br> O:Beutelsbach Weihersdorf |
| **Match representative postal code** <br> **rcMatchRepresentativePostalCode** | ✓ | ✗ | ✗ | ✗ | CN, SepN | The representative postal code matches to the input | I:76000 <br> O:76*** <br> I:76000 <br> O:76229, RepPostcode: 76*** | I:76000 <br> O:75000 |
| **Match length** <br> **rcMatchLength** | ✓ | ✓ | ✓ | ✗ | CN, SN, SepN | The length of the input and the result is the same | I:Beutelsbach - Weihersdorf <br> O: Beutelsbach Weihersdorf <br> I: Karlsruhe <br> O: Stuttgart | I:Beu Wei <br> O: Beutelsbach Weihersdorf |
| **Match exact** <br> **rcMatchExact** | ✗ | ✓ | ✓ | ✗ | CN, SN, SepN, AN | Result is *Multiword Prefix* of input and Input is *Multiword Prefix* of result | I: Weihersdorf Beutelsbach <br> O: Beutelsbach Weihersdorf | I: Weihersdorf Beutelsdorf <br> O: Beutelsbach Weihersdorf |
| **Match exact** <br> **rcMatchExact** | ✓ | ✗ | ✗ | ✗ | CN, SepN | The strings are identical | I: ab3-b5 <br> O: AB3 B5 | I: AB3B5 <br> O: AB3 B5 |
| **Match exact** <br> **rcMatchExact** | ✗ | ✗ | ✗ | ✓ | (none) | The strings are identical | I: 42 <br> O: 42 | I: 42a <br> O: 42 |
| **Match except separators** <br> **rcMatchExceptSeparators** | ✓ | ✓ | ✓ | ✗ | (delegated to Match exact) | After removing all separators Match exact applies | I:Media Park <br> O:Mediapark | I:Media Par <br> O:Mediapark |
| **Match quality high** <br> **rcMatchQualityHigh** | ✓ | ✓ | ✓ | ✗ | (various) | The rating algorithm of the gpGeocoder indicates a high quality | I: Stampfstr <br> O: Stumpfstr | I: Bbuch <br> O: Beutelsbach Weihersdorf |
| **Match quality medium** <br> **rcMatchQualityMedium** | ✓ | ✓ | ✓ | ✗ | (various) | The rating algorithm of the gpGeocoder indicates a medium quality | I: Stamfstr <br> O: Stumpfstr | I: Karlsruhe <br> O: Berlin |
| **Match quality low** <br> **rcMatchQualityLow** | ✓ | ✓ | ✓ | ✓ | (various) | The rating algorithm of the gpGeocoder indicates a low quality | I: Karlsruhe <br> O: Berlin | I: Karlsruhe <br> I: Stamfstr <br> O: Stumpfstr |
| **CityFieldPermutation** <br> **rcCityFieldPermutation** | ✗ | ✓ | ✗ | ✗ | CN, SN, SepN, AN | Not(Considering input City field and result City field Input is Multiword Prefix | I: C: Durlach <br> O:C: Karlsruhe <br> C2: Durlach | I: C: <br> C2: Durlach <br> O: C: Karlsruhe <br> C2: Durlach |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | of result is fulfilled. The same condition must hold for city2.) | | |
| **Match exact except additions**<br><br>**rcMatchExactExceptAdditions** | ⛔ | ⛔ | ⛔ | ✅ | (delegated to Match exact) | If additions are removed, input and result Match exact | I: 42a<br>O: 42 | I: 4 2a<br>O: 42 |

**Notes**

> The examples in above tables are only for explanatory purposes. The authoritative information is given in the definition.
> Note that always exactly one of Match quality high, Match quality medium or Match quality low is true.

*PTV Geoplatform, created on Tue Sep 7 15:37:22 2010*