



Internationalisierung bei der Software-Entwicklung

Thomas Tampe

Leiter Solution Center
Logics Software GmbH
München

Code once, run anywhere ...

... in jeder Sprache und in jedem Land ...

- in Deutsch oder Französisch ...
- in England oder den USA ...
- in Türkisch oder Polnisch ...
- in Griechenland oder Bulgarien ...
- in Hebräisch oder Farsi ...
- in Japan oder China ...

Agenda

- Begriffe und Definitionen
- Zeichensätze – Historie und Grundlagen
- Internationalisierung und Lokalisierung
- Anwendungsarchitekturen
 - Checkliste
 - Werkzeuge
 - Beispiel
- Quellen und Links

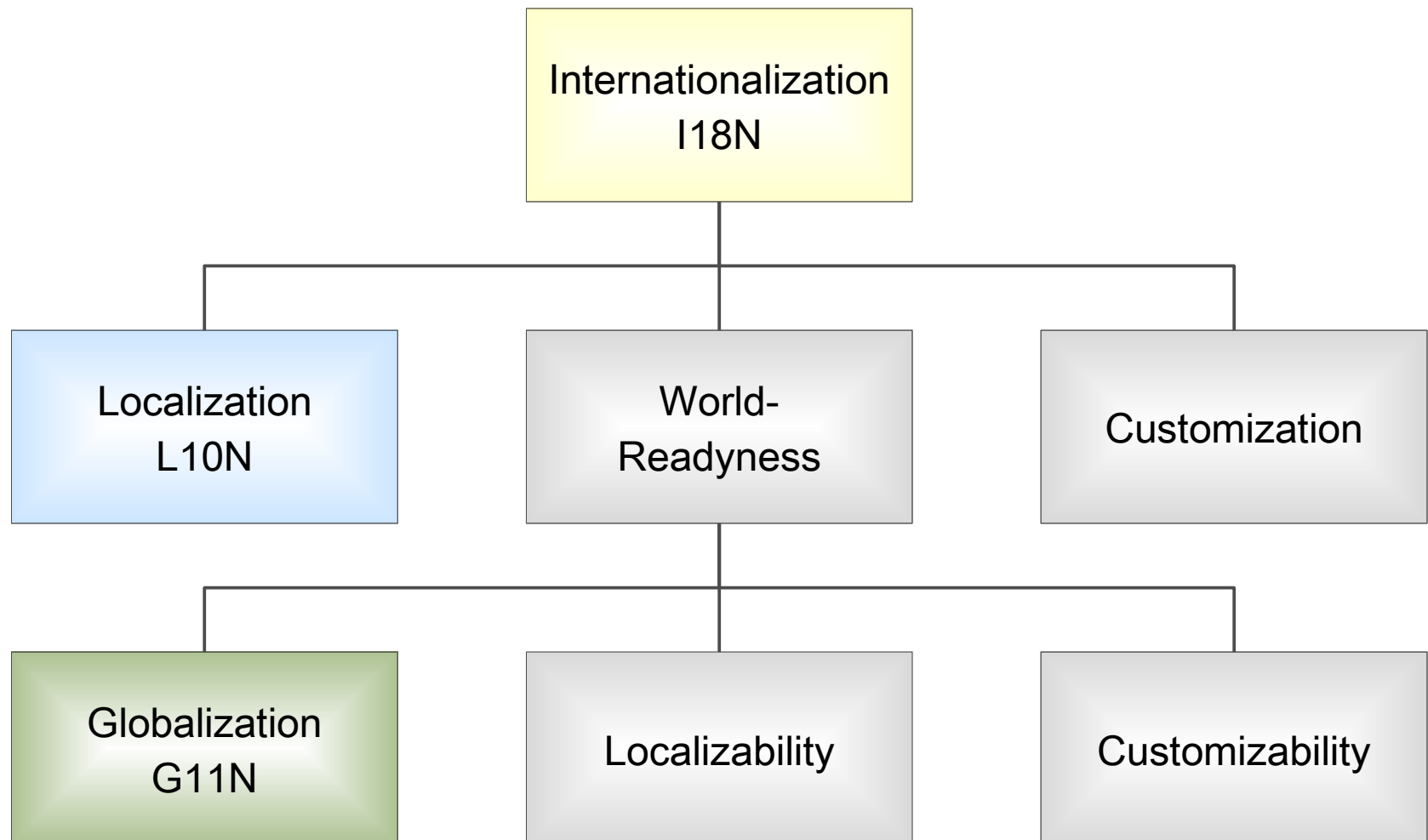
Agenda

- **Begriffe und Definitionen**
- Zeichensätze – Historie und Grundlagen
- Internationalisierung und Lokalisierung
- **Anwendungsarchitekturen**
 - Checkliste
 - Werkzeuge
 - Beispiel
- **Quellen und Links**

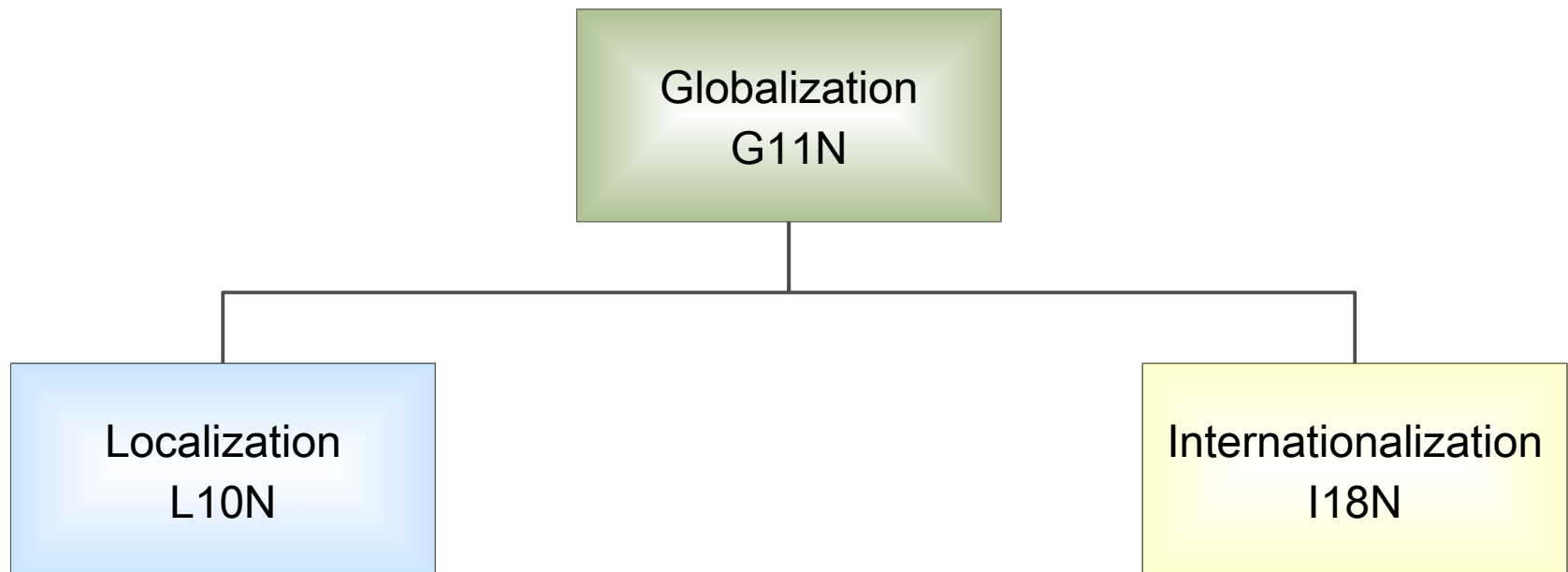
Begriffe und Definitionen

- Internationalisierung
 - I18N = Internationalization
 - Grundlage der Lokalisierung
- Lokalisierung
 - L10N = Localization
 - Anpassung an individuelle Sprache und Region
- Globalisierung
 - G11N = Globalization

Microsoft ist anders ...



... als alle anderen



Agenda

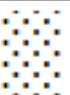
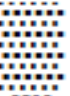

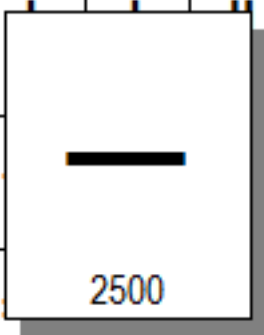
- Begriffe und Definitionen
- Zeichensätze – Historie und Grundlagen
- Internationalisierung und Lokalisierung
- Anwendungsarchitekturen
 - Checkliste
 - Werkzeuge
 - Beispiel
- Quellen und Links

Zeichensätze – Historie


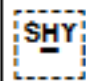
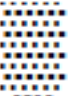

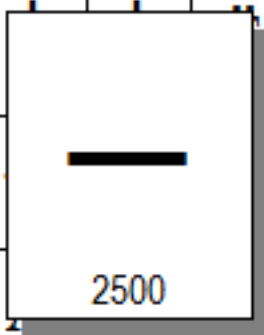
- ASCII (7 Bit)
 - Nur für englische Texte geeignet
- Ein Byte = Ein Zeichen
 - OEM Codepages
 - ANSI Standards
 - ISO-8859-X
- Unicode
 - Code-Points (U+0000 ... U+10FFFF)
 - Oft eingeschränkt auf U+0000 ... U+FFFF

OEM Codepages

437 (OEM – Vereinigte Staaten)

Extended Characters							
008	009	00A	00B	00C	00D	00E	00F
Ç 00C7	É 00C9	á 00E1	 2591	Ł 2514	⌌ 2568	α 03B1	≡ 2261
ü 00FC	æ 00E6	í 00ED	 2592	Ł 2534	⌌ 2564	β 00DF	± 00B1
é 00E9	Æ 00C6	ó 00F3	 2593	⌌ 252C	⌌ 2565	Γ 0393	≥ 2265
â 00E2	ô 00F4	ú 00FA	 2500	⌌ 253C	⌌ 2552	π 03C0	≤ 2264
ä 00E4	ö 00F6	ñ 00F1		Σ 03A3	∫ 2320		
à 00E0	ò 00F2	Ñ 00D1	2561	253C	2552	σ 03C3	Ј 2321
å 00E5	û 00FB	ª 00AA	⌌ 2562	⌌ 255E	π 2553	μ 00B5	÷ 00F7

850 (OEM – Multilingual Lateinisch 1)

Extended Characters							
008	009	00A	00B	00C	00D	00E	00F
Ç 00C7	É 00C9	á 00E1	 2591	Ł 2514	δ 00F0	Ó 00D3	 00AD
ü 00FC	æ 00E6	í 00ED	 2592	Ł 2534	Đ 00D0	β 00DF	± 00B1
é 00E9	Æ 00C6	ó 00F3	 2593	⌌ 252C	Ê 00CA	Ô 00D4	≡ 2017
â 00E2	ô 00F4	ú 00FA	 2500	⌌ 253C	⌌ 2552	Ò 00D2	¾ 00BE
ä 00E4	ö 00F6	ñ 00F1		Σ 03A3	¶ 2320		
à 00E0	ò 00F2	Ñ 00D1	2561	253C	0131	Õ 00D5	§ 00A7
å 00E5	û 00FB	ª 00AA	⌌ 2562	⌌ 255E	π 2553	μ 00B5	÷ 00F7

ANSI Standards

1252 (ANSI – Lateinisch I)

Extended Characters							
008	009	00A	00B	00C	00D	00E	00F
€ 204C	␣ 0090	␣ 00A0	◊ 00B0	À 00C0	Ð 00D0	à 00E0	ð 00F0
␣ 00B1	‘ 2018	ì 00A1	± 00B1	Á 00C1	Ñ 00D1	á 00E1	ñ 00F1
‚ 201A	’ 2019	ç 00A2	² 00B2	Â 00C2	Ò 00D2	â 00E2	ò 00F2
f 0192	“ 201C	£ 00A3	~ 00B3	Ä 00C4	Ó 00D3	ã 00E3	ó 00F3
” 201E	” 201D	⊘ 00A4	~ 00B4	ä 00C4	ô 00D4	~ 00E4	~ 00F4
... 2026	• 2022	¥ 00A5	~ 00B5	å 00C5	õ 00D5	~ 00E5	~ 00F5
† 2020	— 2013	¦ 00A6	¶ 00B6	Æ 00C6	Ö 00D6	æ 00E6	ö 00F6

1251 (ANSI – Kyrillisch)

Extended Characters							
008	009	00A	00B	00C	00D	00E	00F
Ђ 0402	ђ 0452	␣ 00A0	◊ 00B0	А 0410	Р 0420	а 0430	р 0440
Ѓ 0403	‘ 2018	Ў 040E	± 00B1	Б 0411	С 0421	б 0431	с 0441
‚ 201A	’ 2019	Ў 045E	І 0406	В 0412	Т 0422	в 0432	т 0442
Ѓ 0453	“ 201C	Ј 0408	~ 00B3	Д 0414	У 0423	Г 0433	у 0443
” 201E	” 201D	⊘ 00A4	~ 00B4	д 0414	Ф 0424	д 0434	ф 0444
... 2026	• 2022	Г 0490	~ 00B5	е 0415	Х 0425	е 0435	х 0445
† 2020	— 2013	¦ 00A6	¶ 00B6	Ж 0416	Ц 0426	ж 0436	ц 0446

ISO-8859-X

ISO 8859-1 (Lateinisch I)

Extended Characters							
008	009	00A	00B	00C	00D	00E	00F
XXX 0080	DCS 0090	NB SP 00A0	◊ 00B0	À 00C0	Ð 00D0	à 00E0	ð 00F0
XXX 0081	PU1 0091	ı 00A1	± 00B1	Á 00C1	Ñ 00D1	á 00E1	ñ 00F1
BPH 0082	PU2 0092	ç 00A2	² 00B2	Â 00C2	Ò 00D2	â 00E2	ò 00F2
NBH 0083	STS 0093	£ 00A3	³ 00B3	Ä 00C4	Ó 00D3	ä 00E3	ó 00F3
IND 0084	CCH 0094	◊ 00A4	⁴ 00B4	Å 00C5	Ô 00D4	å 00E4	ô 00F4
NEL 0085	MW 0095	¥ 00A5	⁵ 00B5	Ä 00C5	Õ 00D5	ä 00E5	õ 00F5
SSA 0086	SPA 0096	 00A6	¶ 00B6	Æ 00C6	Ö 00D6	æ 00E6	ö 00F6

ISO 8859-7 (Griechisch)

Extended Characters							
008	009	00A	00B	00C	00D	00E	00F
XXX 0080	DCS 0090	NB SP 00A0	◊ 00B0	ı 0390	Π 03A0	Û 03B0	π 03C0
XXX 0081	PU1 0091	‘ 02B0	± 00B1	Α 0391	Ρ 03A1	α 03B1	ρ 03C1
BPH 0082	PU2 0092	, 02B2	² 00B2	Β 0392	◊ 03A2	β 03B2	ς 03C2
NBH 0083	STS 0093	£ 00A3	³ 00B3	Γ 0393	◊ 03A3	γ 03B3	σ 03C3
IND 0084	CCH 0094	◊ 00A4	⁴ 00B4	Δ 0394	◊ 03A4	δ 03B4	τ 03C4
NEL 0085	MW 0095	◊ 00A5	⁵ 00B5	◊ 0395	◊ 03A5	ε 03B5	υ 03C5
SSA 0086	SPA 0096	 00A6	¶ 00B6	Α 0396	Ζ 03A6	Φ 03B6	ζ 03C6

Unicode

Control		Latin-1 Supplement					
008	009	00A	00B	00C	00D	00E	00F
XXX	DCS	NB SP	°	À	Ð	à	ð
0080	0090	00A0	00B0	00C0	00D0	00E0	00F0
XXX	PU1	¡	±	Á	Ñ	á	ñ
0081	0091	00A1	00B1	00C1	00D1	00E1	00F1
BPH	PU2	¢	²	Â	Ò	â	ò
0082	0092	00A2	00B2	00C2	00D2	00E2	00F2
NBH	STS	£	³	Ã	Ó	ã	ó
0083	0093	00A3	00B3	00C3	00D3	00E3	00F3
IND	CCH	¤	´	Ä	Ô	ä	ô
0084	0094	00A4	00B4	00C4	00D4	00E4	00F4
NEL	MW	¥	µ	Å	Õ	å	õ
0085	0095	00A5	00B5	00C5	00D5	00E5	00F5
SSA	SPA		¶	Æ	Ö	æ	ö
0086	0096	00A6	00B6	00C6	00D6	00E6	00F6

Cyrillic							
040	041	042	043	044	045	046	047
È	А	Р	а	р	è	Ɔ	Ψ
0400	0410	0420	0430	0440	0450	0460	0470
Ë	Б	С	б	с	ë	Ƶ	ψ
0401	0411	0421	0431	0441	0451	0461	0471
Ђ	В	Т	в	т	ђ	Ѣ	Θ
0402	0412	0422	0432	0442	0452	0462	0472
Ѓ	Г	У	г	у	ѓ	Ѥ	ϑ
0403	0413	0423	0433	0443	0453	0463	0473
Є	Д	Ф	є	ф	ё	ІЄ	V
0404	0414	0424	0434	0444	0454	0464	0474
Ѕ	Е	Х	ѕ	е	х	Є	Ѳ
0405	0415	0425	0435	0445	0455	0465	0475
І	Ж	Ц	і	ж	ц	і	Ѱ
0406	0416	0426	0436	0446	0456	0466	0476

Grundlagen

Darstellung	Zeichensatz	Wert	Bedeutung
—	437 OEM	0xC4	Blockgrafik – einfache horizontale Linie
Д	1251 ANSI	0xC4	Kyrillischer Großbuchstabe DE
Δ	ISO-8859-7	0xC4	Griechischer Großbuchstabe Delta
Ä	ISO-8859-1	0xC4	Lateinischer Großbuchstabe A Umlaut
A	Unicode	U+0040	Lateinischer Großbuchstabe A
Α	Unicode	U+0391	Griechischer Großbuchstabe Alpha
А	Unicode	U+0410	Kyrillischer Großbuchstabe A

- Zeichensatz und Wert legen Bedeutung fest
- Codierung dient Speicherung und Transport
- Font (Schriftart) stellt Zeichen dar
- There Ain't No Such Thing As Plain Text. (Joel Spolsky)

Beispiel – Codierung

- Erwartete Textdarstellung

l'élection présidentielle française

- Text codiert als UTF-8 aber als ISO-8859-1
angezeigt

lâ€™™ Ã©lection prÃ©sidentielle franÃ§aise

- Text codiert als ISO-8859-1 aber als UTF-8
angezeigt

l?|?ction pr?|?dentielle fran?|?se

<http://www.alanwood.net/unicode/htmlunicode.html>

Unicode / UTF-8

- **Aktueller Standard**
 - Derzeit $17 * 2^{16} = 1.114.112$ Zeichen
 - Zeichen U+0000 ... U+00FF identisch zu ISO-8859-1
- **Codierung UTF-8**
 - maximal 2^{31} Zeichen
 - Variable Länge (1 – 6 Byte pro Zeichen)

Zeichen	Buchstabe y	Buchstabe Ä	Zeichen ©	Eurozeichen €
Unicode	U+0079	U+00C4	U+00A9	U+20AC
Unicode binär	00000000 01111001	00000000 11000100	00000000 10101001	00100000 10101100
UTF-8 binär	01111001	11000011 10000100	11000010 10101001	11100010 10000010 10101100
UTF-8	0x79	0xC3 0x84	0xC2 0xA9	0xE2 0x82 0xAC

- **Kompatibel mit 7-Bit ASCII**
 - Codierung als einzelnes Byte
 - NUL (0x00) markiert das Ende einer Zeichenreihe
- **Überprüfbare Codierung**
 - Teilsequenzen können erkannt werden
 - Anfang einer Sequenz ist ermittelbar
- **Unterstützung in modernen Betriebssystemen**
 - Windows (NT, 2000, XP, Vista, 7)
 - Unix (Solaris, AIX, HP-UX, ...)
 - Linux (Debian, Ubuntu, SUSE, ...)

Agenda

- Begriffe und Definitionen
- Zeichensätze – Historie und Grundlagen
- **Internationalisierung und Lokalisierung**
- **Anwendungsarchitekturen**
 - Checkliste
 - Werkzeuge
 - Beispiel
- Quellen und Links

Internationalisierung

- Intern Unicode (U+0000 ... U+FFFF)
 - Transport und Speicherung als UTF-8
- Texte zur Laufzeit laden
 - Verschiedene Pluralformen
 - Platzhalter für Werte
 - Anzeigeformate für Datum / Uhrzeit, Zahlen, ...
- Medien zur Laufzeit laden
 - Symbole
 - Audio- und Video-Clips

Internationalisierung

- Eingabemethoden
 - Tastaturkürzel für Menübefehle
 - Unabhängigkeit vom Tastaturlayout
- Layout der Benutzeroberfläche
 - Unterschiedliche Länge der Texte
 - Schreibrichtung (Bidirektionale Texte)
- Druckausgaben
 - Installierte Druckertreiber
 - Unterstützte Papierformate

Internationalisierung

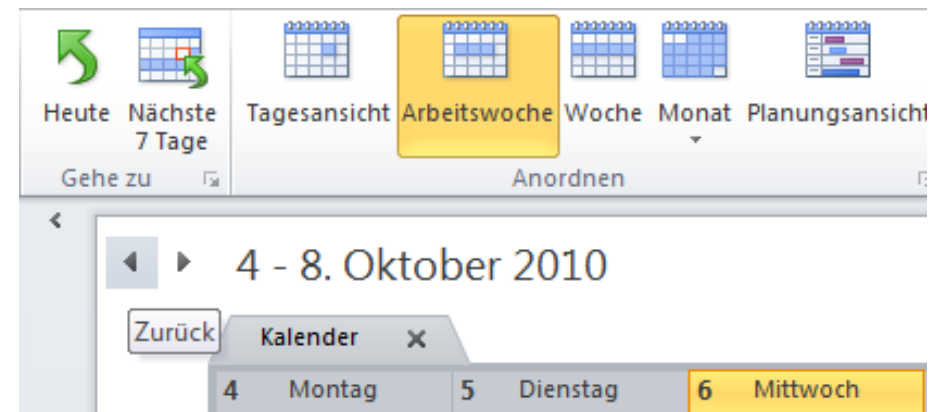
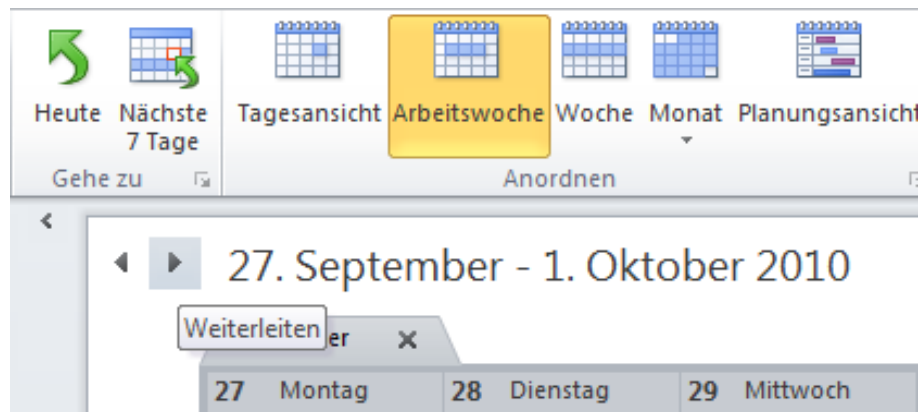
- Interpretation der Eingabedaten
 - Gibt es auch andere Ziffern als 0 ... 9 ?
 - Was sind Großbuchstaben / Kleinbuchstaben ?
 - Wo ist das Ä einzusortieren ?
 - Ist $\frac{1}{4}$ das gleiche wie $\frac{1}{4}$?
- Bedingte Übersetzung für grundlegende Unterschiede eventuell sinnvoll
 - bidirektionale Sprachen (Hebräisch, Arabisch, ...)
 - ostasiatische Sprachen (Chinesisch, Japanisch, ...)

Lokalisierung

- Für jede Locale (Culture) separat
- Übersetzen der (externen) Texte
- Mediale Inhalte bereitstellen
- Falls nicht im Laufzeitsystem verfügbar
 - Formate für Zeit / Datum, Zahlen, ...
 - Währungssymbole
 - Umrechnungsfunktionen für Maßeinheiten
 - Algorithmen für Sortierung
 - Umwandlung Großbuchstaben / Kleinbuchstaben

Lokalisierung

- Setzt sorgfältige Internationalisierung voraus
 - Beispiel aus Microsoft Outlook 2010



- Kontextsensitive Bezeichnungen verwenden
- Konsistente Datumsformate definieren

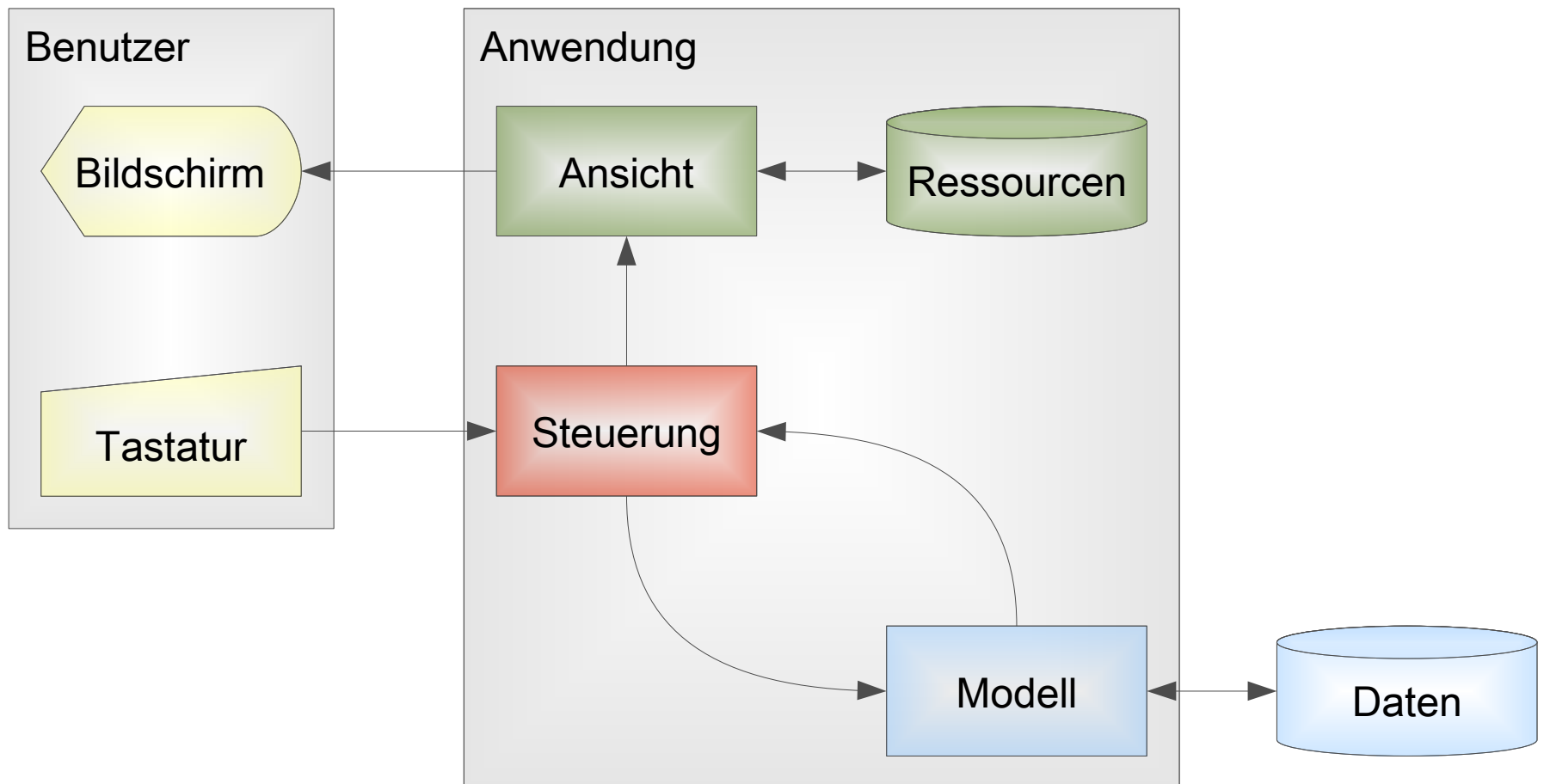
Agenda

- Begriffe und Definitionen
- Zeichensätze – Historie und Grundlagen
- Internationalisierung und Lokalisierung
- **Anwendungsarchitekturen**
 - Checkliste
 - Werkzeuge
 - Beispiel
- Quellen und Links

Model – View – Controller (MVC)

- Trennt Darstellung von Anwendungslogik
- Controller
 - Koordiniert Model und View
- Model
 - Benötigt keine Information zu Sprache / Region
 - Verarbeitet abstrakte Daten
- View
 - Kennt Sprache und Region des Anwenders
 - Stellt abstrakte Daten dar

Desktop Anwendung



Datenbankschnittstelle

- Einzige externe Schnittstelle
- Zugriff über Persistenz-Frameworks
 - TopLink, Hibernate (Java)
 - ADO.NET Entity Framework, NHibernate (.NET)
- Automatische Konvertierung in Objekte
- Entwurfswerkzeuge
 - NetBeans, Eclipse (Java)
 - Visual Studio 2010 (.NET)
- Kein Thema für Internationalisierung

Desktop Anwendung

- Sprache und Region aus Benutzerumgebung
 - Locale (Java) oder Culture (.NET)
 - de_DE ist Deutsch (Deutschland)
 - es_US wäre Spanisch (Vereinigte Staaten)
 - Sprache wichtiger als Region
- Tastaturbelegung oft umschaltbar
 - Symbolische Tastennamen verwenden
 - mögliche Tastenkombinationen berücksichtigen

Desktop Anwendung – Checkliste

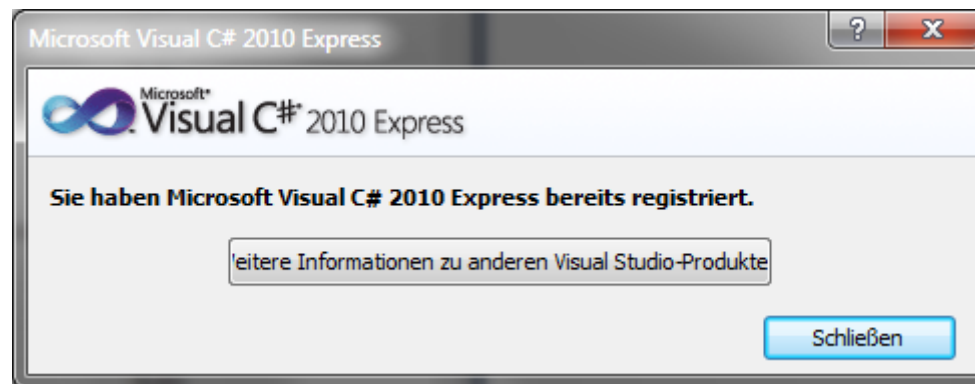
- Fenster
 - Titelzeile
- Menüs
 - Titel und Einträge
 - Tastaturkürzel
- Symbolleisten
 - Symbole
 - Texte für Info-Fenster
- Statuszeile
- Dialogfelder
 - Meldungstexte
- Formate
 - Zeit / Datum
 - Zahlen, Währung
- Kontrollelemente
 - Beschriftungen
 - Inhalte
- Klänge und Bilder

Desktop Anwendung – Werkzeuge

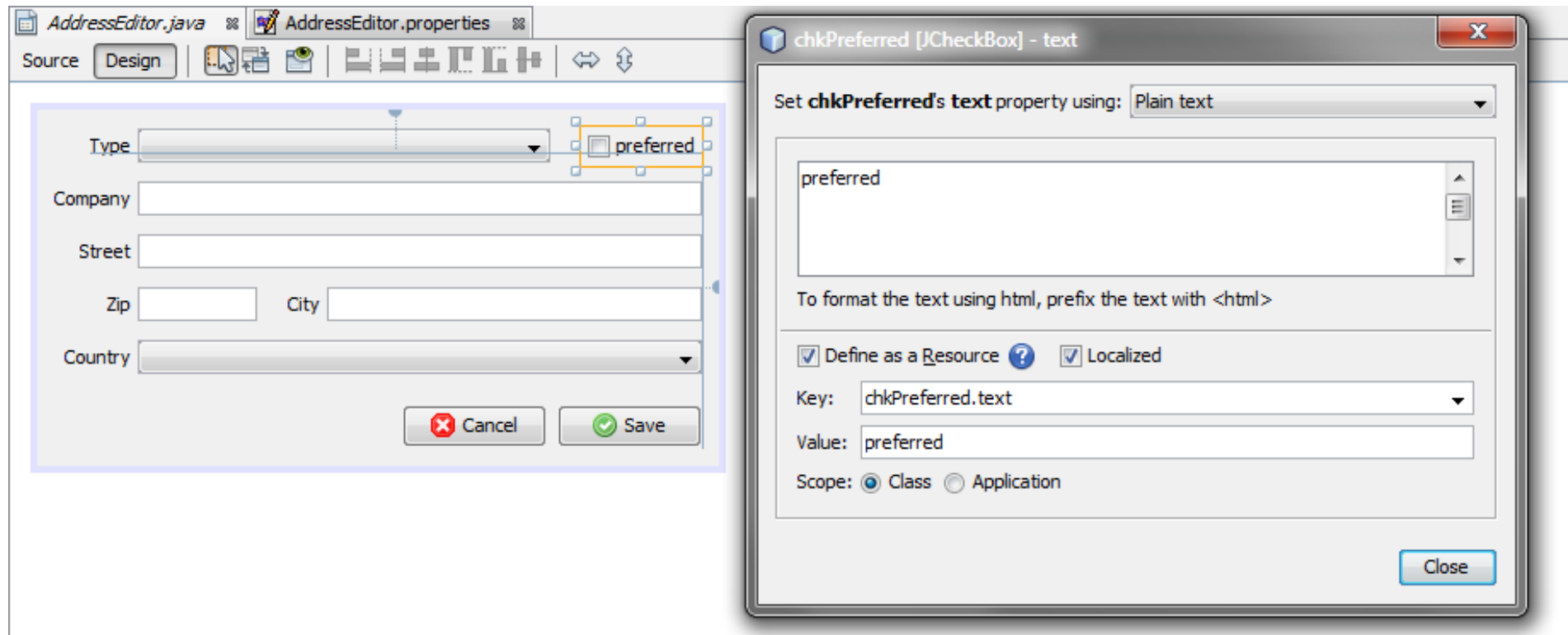
- Integrierte Entwicklungsumgebungen
 - NetBeans, Eclipse (Java)
 - Visual Studio (.NET)
 - Ressourcen automatisch beim Entwurf anlegen
 - *.properties oder ResourceBundle (Java)
 - *.resx, *.resource (.NET)
- Resource Editoren
 - Windows Forms Resource Editor (.NET)

Desktop Anwendung – Layouts

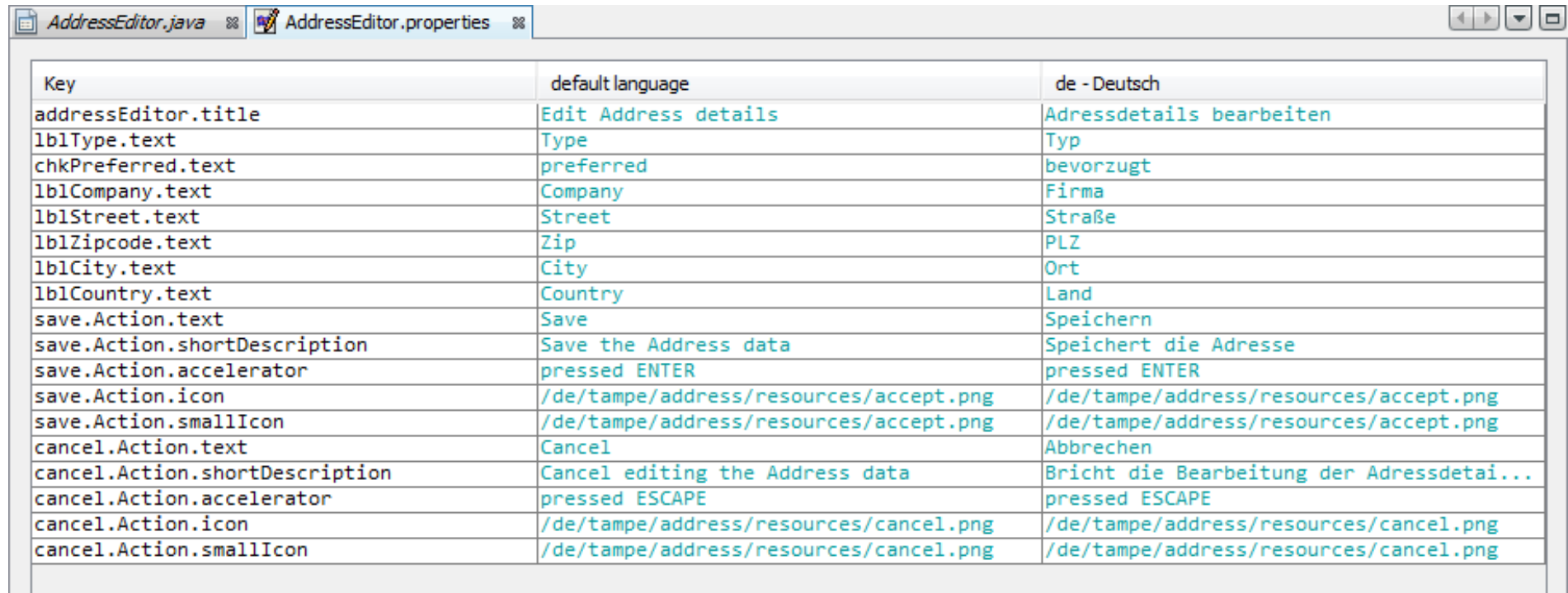
- Flexible Layouts
 - Größe zur Laufzeit an Inhalt anpassen
- Feste Layouts
 - Größe zum Übersetzungszeitpunkt an Inhalt anpassen



Desktop Anwendung – NetBeans

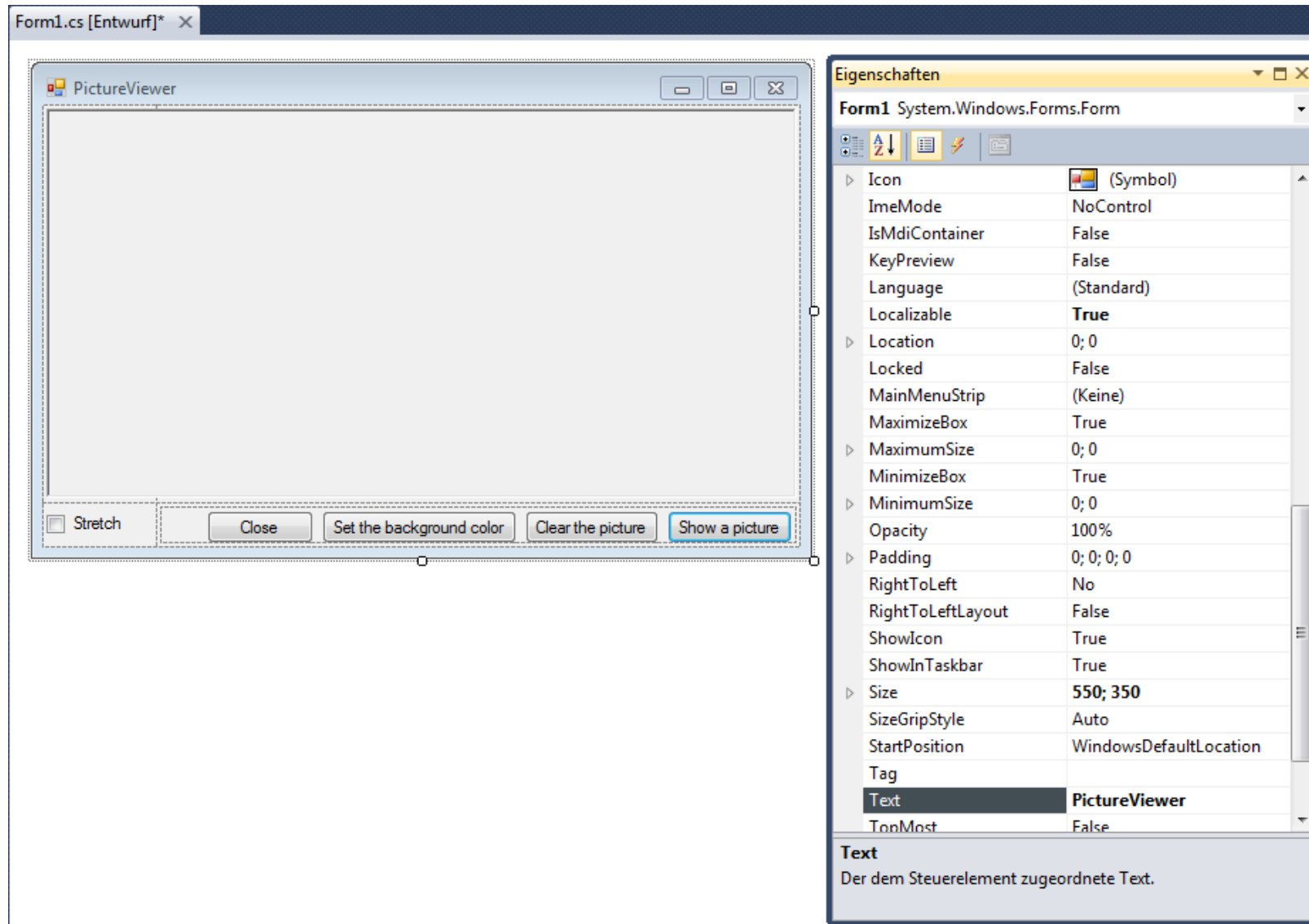


Desktop Anwendung – NetBeans

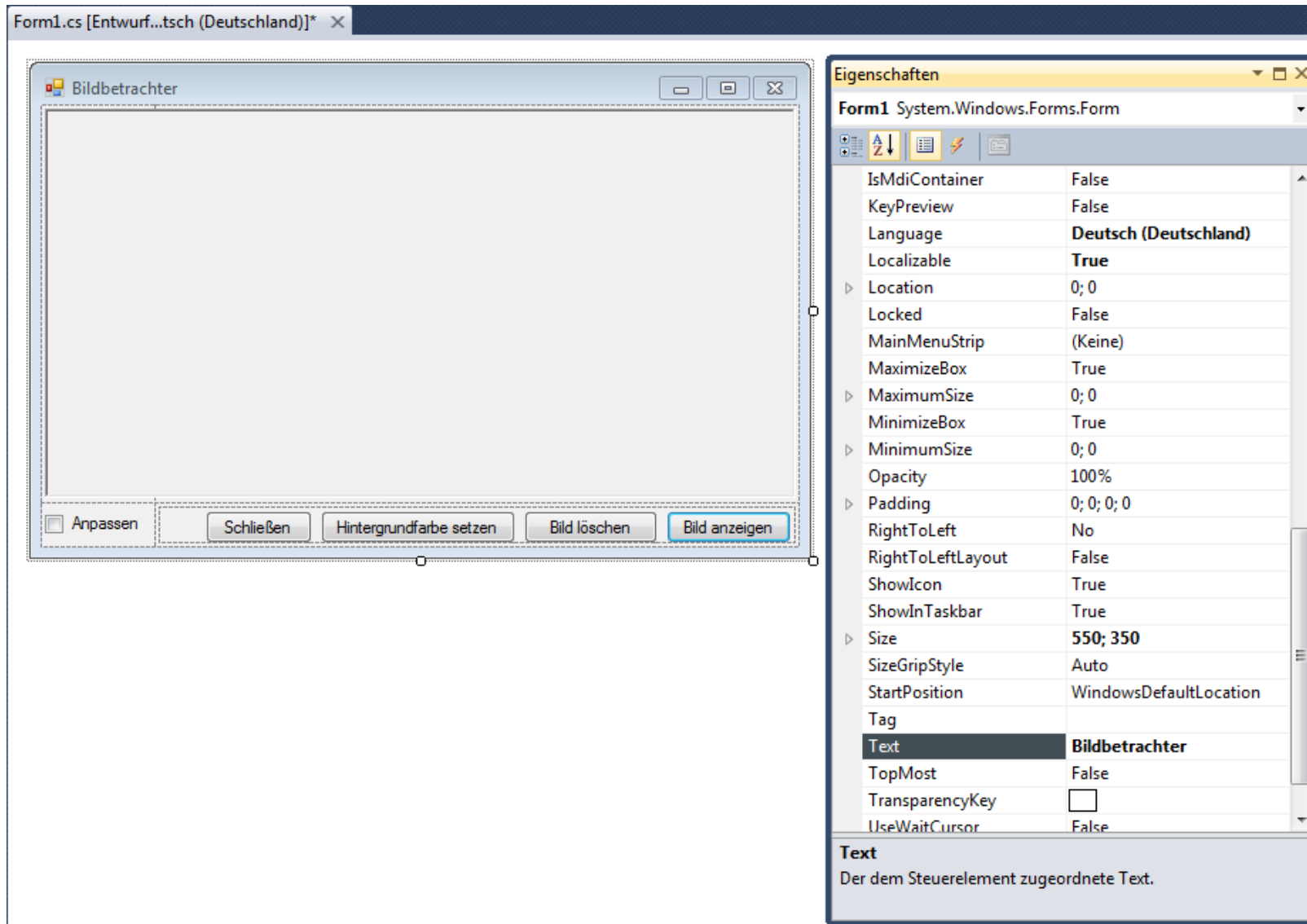


Key	default language	de - Deutsch
addressEditor.title	Edit Address details	Adressdetails bearbeiten
lblType.text	Type	Typ
chkPreferred.text	preferred	bevorzugt
lblCompany.text	Company	Firma
lblStreet.text	Street	Straße
lblZipcode.text	Zip	PLZ
lblCity.text	City	Ort
lblCountry.text	Country	Land
save.Action.text	Save	Speichern
save.Action.shortDescription	Save the Address data	Speichert die Adresse
save.Action.accelerator	pressed ENTER	pressed ENTER
save.Action.icon	/de/tampe/address/resources/accept.png	/de/tampe/address/resources/accept.png
save.Action.smallIcon	/de/tampe/address/resources/accept.png	/de/tampe/address/resources/accept.png
cancel.Action.text	Cancel	Abbrechen
cancel.Action.shortDescription	Cancel editing the Address data	Bricht die Bearbeitung der Adressdetai...
cancel.Action.accelerator	pressed ESCAPE	pressed ESCAPE
cancel.Action.icon	/de/tampe/address/resources/cancel.png	/de/tampe/address/resources/cancel.png
cancel.Action.smallIcon	/de/tampe/address/resources/cancel.png	/de/tampe/address/resources/cancel.png

Desktop Anwendung - Visual Studio



Desktop Anwendung - Visual Studio



Beispiel – Markieren und Kopieren

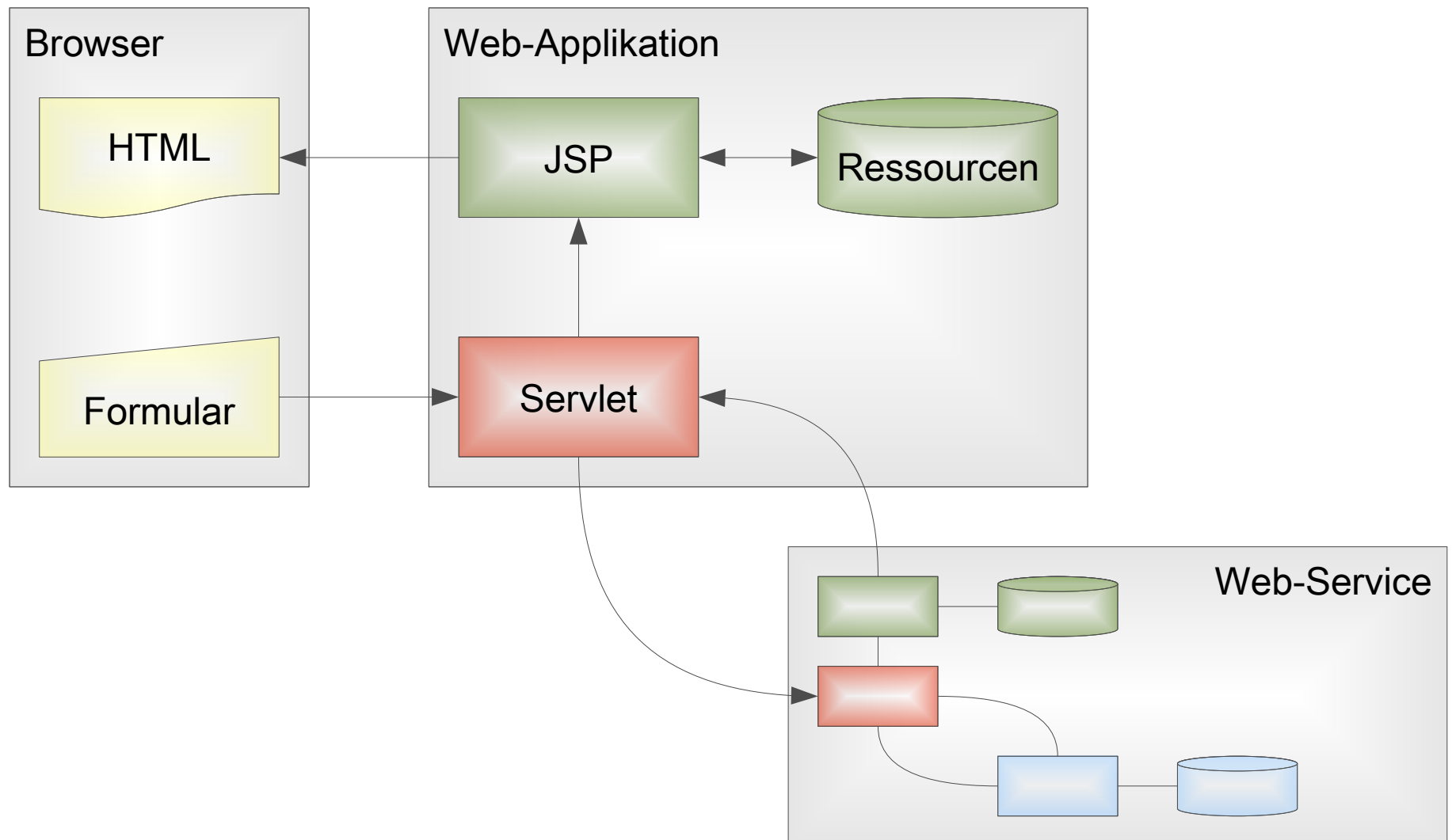
Normaler Text	Buchstaben und Ziffern LTR
Zeichenreihe	ABCD 1234 EFGH
Darstellung	AB <u>CD</u> <u>1234</u> <u>EFGH</u>
Zeichenreihe	AB <u>CD</u> <u>1234</u> <u>EFGH</u>
Zu kopierender Text	CD 1234 EF

Bidirektionaler Text	Buchstaben RTL, Ziffern LTR
Zeichenreihe	ABCD 1234 EFGH
Darstellung	DC <u>BA</u> <u>1234</u> <u>HGFE</u>
Zeichenreihe	<u>ABCD</u> <u>1234</u> <u>EFGH</u>
Zu kopierender Text	AB 1234 GH

Agenda

- Begriffe und Definitionen
- Zeichensätze – Historie und Grundlagen
- Internationalisierung und Lokalisierung
- **Anwendungsarchitekturen**
 - Checkliste
 - Werkzeuge
 - Beispiel
- Quellen und Links

Verteilte Anwendung



HTTP – Schnittstellen

- Client teilt Codierung mit
 - Accept-Charset: utf-8
 - `<form accept-charset="utf-8" ...>`
- Server teilt Codierung und Typ mit
 - HTML-Dokumente
 - Content-Type: text/html; charset=utf-8
 - `<meta http-equiv="content-type" content="text/html; charset=utf-8">`
 - XML-Dokumente
 - Content-Type: application/xml; charset=utf-8
 - `<?xml version="1.0" encoding="UTF-8" ...>`

Verteilte Anwendung

- **Bevorzugte Sprache und Region aus HTTP-Header oder Parameter**
 - Accept-Language: de-DE, en-US;q=0.5
 - ... &lang=de-DE ...
- **Service-Aufruf aus URL und / oder Parametern**
 - `http://server.domain/controller/service.jsp?param1=value1& ...`
- **Service-spezifische JSP / ASP aus Sprache und Region**

Verteilte Anwendung – Checkliste

- HTML-Seiten
 - Texte
- Stylesheets
 - Symbole
- Scripts
 - Meldungstexte
 - Tastaturkürzel
- Klänge
- Bilder
- Ressourcen
 - Texte
 - Zeit- / Datums- und Zahlformate
- Frameworks

Verteilte Anwendung – Werkzeuge

- Integrierte Entwicklungsumgebungen
 - NetBeans, Eclipse, Visual Studio
- Editor für HTML / JSP / ASP
 - DreamWeaver, Kompozer, ...
 - Nur sichtbare Teile der Seite übersetzen
- Programm zum Dateivergleich
 - ExamDiff, WinMerge, ...
 - Strukturänderungen der Seite übernehmen
- Windows Forms Resource Editor (.NET)

Beispiel – Online-Shop

- Bestellung per Internet auch aus dem Ausland
- Rechnung mit ausgewiesener Umsatzsteuer

- Model liefert nur Netto-Betrag
- View berechnet Steuer- und Brutto-Betrag

- Controller liefert Sprache und Region
- Model liefert Netto-, Steuer- und Brutto-Betrag
- Land und Steuersatz gehören zum Model

- Rechnung wird von View erstellt
 - Text, Datumsformate, Währungsformate, ...

Agenda

- Begriffe und Definitionen
- Zeichensätze – Historie und Grundlagen
- Internationalisierung und Lokalisierung
- Anwendungsarchitekturen
 - Checkliste
 - Werkzeuge
 - Beispiel
- **Quellen und Links**

Quellen und Links

Unicode and Character Sets

<http://www.joelonsoftware.com/articles/Unicode.html>

UTF-8, a transformation format of ISO 10646

<http://www.ietf.org/rfc/rfc3629.txt>

Unibook™ Character Browser

<http://unicode.org/unibook/>

Encoding and Localization

[http://msdn.microsoft.com/en-us/library/h6270d0z\(v=VS.100\).aspx](http://msdn.microsoft.com/en-us/library/h6270d0z(v=VS.100).aspx)

Microsoft .NET Internationalization

<http://msdn.microsoft.com/en-us/goglobal/bb688096.aspx>

ASP.NET Wiki: Internationalization

<http://wiki.asp.net/page.aspx/55/internationalization/>

ADO.NET Entity Framework

<http://msdn.microsoft.com/en-us/library/bb399572.aspx>

Open Source Persistence Frameworks in C#

<http://csharp-source.net/open-source/persistence>

Robust internationalization with GTK+

<http://www.ibm.com/developerworks/aix/library/au-internatlgtk/index.html>

Create Great-Looking GUIs With NetBeans IDE 5.5

http://java.sun.com/developer/technicalArticles/tools/nb_guibuilder/



An old joke in the internationalization community goes like this:

A person who speaks three languages is called trilingual.

And a person who speaks two languages is called bilingual.

So what do you call someone who only speaks one language?

American.

<http://www.eclipse.org/articles/Article-Internationalization/how2I18n.html>



Thomas Tampe

thomas.tampe@arcor.de
+49 (172) 9 57 17 05

t.tampe@logics.de
+49 (89) 55 24 04 32