



Business Case Model for use of PDQ in Data Capture QA

A White Paper by K. Bradley Paxton, Ph.D.

15 Jun 2011

Executive Summary

In doing forms data capture, whether with just human data entry keying from paper forms or with Optical Character Recognition (OCR), Optical Mark Recognition (OMR) and Key From Image (KFI), it has been customary to employ manual methods for data quality assurance. These methods involve a process we refer to as Double Key & Verify (DK&V), wherein one keyer is asked to key a particular data field, and then another keyer is asked to key the same field (preferably without collusion). If the results from both keyers agree with the sampled field, then the sampled data field is deemed to be correct. If they do not, then a third party is usually employed to divine the correct answer. This classic DK&V process is slow and costly; so, in practice, the amount of data sampled for Quality Assurance (QA) purposes is often smaller than desired for useful statistically valid results.

At ADI, we have developed a new automated approach to data capture QA that we call Production Data Quality (PDQ). In brief, the PDQ system employs independent data capture engines to sample the production data capture images and results, and to quickly and cost-effectively determine the truth of the sampled fields. By cost-effective, we mean that the amount of human keying needed for QA may be reduced by a factor of 40 or more. In addition, when the truth is known, data quality accuracy can be measured precisely, root-cause analysis is enabled, and data capture system improvements made more rapidly.

A version of this PDQ system was developed and used in the U.S. Census Bureau's 2010 Decennial Census and was a great success, providing near-real-time feedback to Census management of data capture issues and assurance that the required data quality metrics were being met.

This short white paper describes a useful business case model to use to estimate the keyer QA savings you can achieve using PDQ. It is not our intent here to explain the internal workings of PDQ in detail, but a process flow chart is shown in Appendix 1. The model's nominal values herein are built around the CMS-1500 Health Insurance Claim Form, examples of which are shown in Appendices 2 and 3, however, the model can be used for any form type you wish by simply changing the model inputs appropriately.

How to Use the Business Case Model

Using the model discussed in this white paper is easy to do. Basically, you just put in values that describe your form and data entry process in the input section; and the savings to you if you were to employ PDQ to do your data entry QA are immediately calculated. If you want to try a different estimate of a particular value, you can easily change it and see “what if.”

Below we will briefly describe the inputs, and then show you an example. Table 1 below lists the ten inputs and some nominal values.

Assumptions (Inputs)	Values
Form Type	CMS-1500
OCR Accept Rate	0.8
Keystrokes/Hour	6000
Burdened Keying Cost (\$/Hour)	15
Form Volume (Megaforms)	1
Characters/Form	1000
OCR QA Sampling Rate	0.01
KFI QA Sampling Rate	0.05
DK&V Factor	2
PDQ Efficiency over DK&V	40

Table 1 – The Ten Inputs to the PDQ Business Case Model
(With some typical values)

Form Type

Here you just record the name of the form type from which you are capturing data, just to label the calculation so you can remember what you did. Here, in our example, we have assumed the form type is the CMS-1500 Health Insurance Claim Form, formerly referred to as “HCFA.”

OCR Accept Rate

The OCR Accept Rate is that fraction of your data capture work being read automatically by your recognition system. Here, we assume 0.8 (or 80%). The remainder of the work, in this case 0.2 (or 20%) is the Reject Rate, and this production work is sent to your human keyers because the OCR is unsure about the answer. If you are doing all human data entry keying and not using automation at all, then simply enter 0.0 for the OCR Accept Rate (equivalent to rejecting everything to keyers).

Keystrokes/Hour

This is just the average number of keystrokes you estimate your keyers punch per hour under normal daily (not peak) working conditions. Here, we have assumed 6,000 keystrokes per hour.

Burdened Keying Cost (\$/Hour)

This one is a little tricky, because it needs to be the total cost of employing a keyer in your enterprise. Ideally, it would include not only the hourly wage, but also the cost of equipment, space, heat and light, etc. You may need your CPA. Here, we have assumed an approximate U.S. minimum wage with a 2X burden rate, or about \$15/hour.

Form Volume (Megaforms)

Here we use a new term, "Megaforms," which as you can probably guess means a volume of one million forms. We introduced this term some time back, as the volume at which data capture automation should seriously be considered. Here, we assume only one "Megaform," which is also easy to mentally scale up.

Characters/Form

This input is the average number of characters placed on a form by the respondent, whether machine print or handprint. We have estimated 1000 characters based on a CMS-1500 form.

OCR QA Sampling Rate

This would be the rate at which you wish to statistically sample the (accepted) output from your OCR system for QA purposes. Here, we have assumed a sampling rate of 0.01 (or 1%), which was what was actually done in Census 2010.

KFI QA Sampling Rate

This is the rate at which you wish to sample the results of your human keying for QA purposes, whether Key From Paper (KFP) or Key From Image (KFI). Here, we assume 0.05 (or 5%), which is what was actually done in Census 2010.

DK&V Factor

This factor just accounts for the extra keying required when Double Key & Verify is employed. We have only assumed a factor of two here, but often it is really about 2.2 in practice due to the "verify" part.

PDQ Efficiency over DK&V

This number can really be from about 30 for very discriminating studies like the Census up to 100 or more, depending on the overall accuracy of your system that you are trying to measure. You can think of it as the reduction factor for keying effort in doing the QA by employing PDQ instead of DK&V. Here, we assume a nominal value of 40 (you may envision your army of QA keyers being divided by 40 if you use PDQ). In order for you to really believe this one, you may have to try it and see how it works with your actual forms and your workflow, but this is where the money is.

The “Live” PDQ Business Case Model

If we start by assuming that all the above assumptions are OK for now, and run the model, then you get the following:

Assumptions (Inputs)	Values
Form Type	CMS-1500
OCR Accept Rate	0.8
Keystrokes/Hour	6000
Burdened Keying Cost (\$/Hour)	15
Form Volume (Megaforms)	1
Characters/Form	1000
OCR QA Sampling Rate	0.01
KFI QA Sampling Rate	0.05
DK&V Factor	2
PDQ Efficiency over DK&V	40

Results (Outputs)	Values
Total Character Volume	1.00E+09
OCR Character Volume	8.00E+08
KFI Character Volume	2.00E+08
QA Character Volume	3.60E+07
QA Keying Time in Hours	6.0E+03
QA Cost using DK&V	\$90,000
QA Cost using PDQ	\$2,250

QA Savings using PDQ	\$87,750
-----------------------------	-----------------

This white paper is just focused on the obvious QA cost savings, but you can also make significant continuous quality improvements using PDQ, which will be the topic of another paper.

For further information, or to get an Excel workbook of this model, contact:

Brad Paxton

ADI, LLC

200 Canal View Boulevard, Suite 100

Rochester, NY 14623

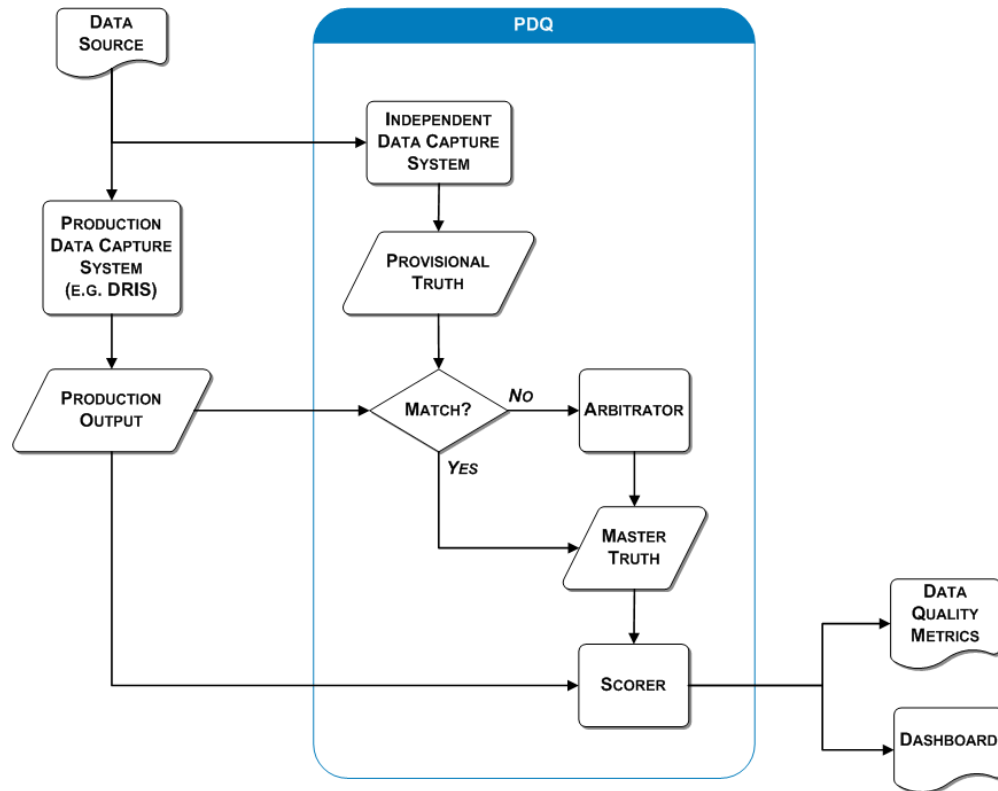
(585) 239-6057

brad.paxton@adillc.net

<http://www.adillc.net>

Appendix 1 – PDQ Process Flow Diagram

We have intentionally not tried to explain the inner workings of PDQ in this note that is focused just on QA cost savings. However, here is a high-level flow chart of PDQ in case you're interested.



Appendix 2 – Example of a Machine-Printed CMS-1500 Form (Synthetic Data)

1500
HEALTH INSURANCE CLAIM FORM

APPROVED BY NATIONAL UNIFORM CLAIM COMMITTEE 08/05

PICA
PICA

1. MEDICARE <input type="checkbox"/> MEDICAID <input type="checkbox"/> TRICARE CHAMPUS (Sponsor's SSN) <input checked="" type="checkbox"/> CHAMPVA <input type="checkbox"/> GROUP HEALTH PLAN (SSN or ID) <input type="checkbox"/> FECA BLX (LUNG) (SSN) <input type="checkbox"/> OTHER <input type="checkbox"/>		1a. INSURED'S I.D. NUMBER (For Program in item 1) 53MB985G3TR744C99STC46VLE7FO	
2. PATIENT'S NAME (Last Name, First Name, Middle Initial) West, Imelda M.		4. INSURED'S NAME (Last Name, First Name, Middle Initial) Richard, Amy M.	
3. PATIENT'S BIRTH DATE MM DD YY 01 15 1976 M <input type="checkbox"/> F <input checked="" type="checkbox"/>		7. INSURED'S ADDRESS (No., Street) 5374 Faye Circle	
5. PATIENT'S ADDRESS (No., Street) 5374 Faye Circle		6. PATIENT RELATIONSHIP TO INSURED Self <input type="checkbox"/> Spouse <input checked="" type="checkbox"/> Child <input type="checkbox"/> Other <input type="checkbox"/>	
8. PATIENT STATUS Single <input checked="" type="checkbox"/> Married <input type="checkbox"/> Other <input type="checkbox"/>		8. PATIENT STATUS Single <input checked="" type="checkbox"/> Married <input type="checkbox"/> Other <input type="checkbox"/>	
9. OTHER INSURED'S NAME (Last Name, First Name, Middle Initial) Terrell, Garrett R.		10. IS PATIENT'S CONDITION RELATED TO: Employed <input checked="" type="checkbox"/> Full-Time Student <input type="checkbox"/> Part-Time Student <input type="checkbox"/>	
11. INSURED'S POLICY GROUP OR FECA NUMBER 89UU768B2CB921D09NNV25ATN9NR		11. INSURED'S POLICY GROUP OR FECA NUMBER 89UU768B2CB921D09NNV25ATN9NR	
a. OTHER INSURED'S POLICY OR GROUP NUMBER 56MC851G3J1875N48RRC03WVWQ5XG		a. EMPLOYMENT? (Current or Previous) <input checked="" type="checkbox"/> YES <input type="checkbox"/> NO	
b. OTHER INSURED'S DATE OF BIRTH MM DD YY 01 28 04 M <input type="checkbox"/> F <input checked="" type="checkbox"/>		b. AUTO ACCIDENT? <input checked="" type="checkbox"/> YES <input type="checkbox"/> NO PLACE (State) NE	
c. EMPLOYER'S NAME OR SCHOOL NAME Mount Saint Vincent University		c. OTHER ACCIDENT? <input type="checkbox"/> YES <input checked="" type="checkbox"/> NO	
d. INSURANCE PLAN NAME OR PROGRAM NAME Sidney Hillman Health Center		10d. RESERVED FOR LOCAL USE	
12. PATIENT'S OR AUTHORIZED PERSON'S SIGNATURE I authorize the release of any medical or other information necessary to process this claim. I also request payment of government benefits either to myself or to the party who accepts assignment below. <i>Imelda West</i> 06/02/10		13. INSURED'S OR AUTHORIZED PERSON'S SIGNATURE I authorize payment of medical benefits to the undersigned physician or supplier for services described below. <i>Amy Richard</i>	
14. DATE OF CURRENT ILLNESS (First symptom) OR INJURY (Accident) OR PREGNANCY (LMP) MM DD YY 02 13 2010		15. IF PATIENT HAS HAD SAME OR SIMILAR ILLNESS GIVE FIRST DATE MM DD YY 8 20 2002	
17. NAME OF REFERRING PROVIDER OR OTHER SOURCE Sidney Hillman Health Center		16. DATES PATIENT UNABLE TO WORK IN CURRENT OCCUPATION FROM MM DD YY TO MM DD YY 4 2 01 6 9 01	
17a. 999		18. HOSPITALIZATION DATES RELATED TO CURRENT SERVICES FROM MM DD YY TO MM DD YY 4 2 01 4 6 01	
17b. NPI 999		20. OUTSIDE LAB? <input checked="" type="checkbox"/> YES <input type="checkbox"/> NO \$ CHARGES 1234 23	
21. DIAGNOSIS OR NATURE OF ILLNESS OR INJURY (Relate Items 1, 2, 3 or 4 to Item 24E by Line) 1. LUZP8, 0346 3. XNL2, 8480		22. MEDICAID RESUBMISSION CODE ORIGINAL REF. NO. VGX75ECD3LD DSH03YTT9DV715FS67	
2. CGT3 5932 4. HAU4 8777		23. PRIOR AUTHORIZATION NUMBER YGE12NFS8WF145ZG8433KRYR63KQT	
24. A. DATE(S) OF SERVICE From MM DD YY To MM DD YY		B. PLACE OF SERVICE	
C. EMG		D. PROCEDURES, SERVICES, OR SUPPLIES (Explain Unusual Circumstances) CPT/HCPCS I MODIFIER	
E. DIAGNOSIS POINTER		F. \$ CHARGES	
G. DAYS OR UNITS		H. EPST Family Plan	
I. ID DUAL		J. RENDERING PROVIDER ID #	
12 14 02 01 15 03 C9 FW 42CQ69 0F 9U 0F Y7 0346 903 84 256 K NPI PRZ83ZGT6LT			
02 21 05 08 21 03 9M QR MG43EZ 4B C1 4B C1 8480 3079 29 124 B NPI YTK50CAJ3EV			
07 16 10 01 10 05 2T QJ 33MY97 3O U0 J5 8B 5932 675 07 6 F NPI QEP70GY2WV			
		NPI	
		NPI	
		NPI	
26. FEDERAL TAX I.D. NUMBER SSN EIN <input checked="" type="checkbox"/>		27. ACCEPT ASSIGNMENT? (or gov't claim, see back) YES <input type="checkbox"/> NO <input checked="" type="checkbox"/>	
28. TOTAL CHARGE \$		29. AMOUNT PAID \$	
30. BALANCE DUE \$		31. BILLING PROVIDER INFO & PH #	
31. SIGNATURE OF PHYSICIAN OR SUPPLIER INCLUDING DEGREES OR CREDENTIALS (I certify that the statements on the reverse apply to this bill and are made a part thereof.) <i>Dr. J. Public</i> 06/18/10		32. SERVICE FACILITY LOCATION INFORMATION 2222 North Lincoln Avenue, York, NE 68467 RR 1 BOX 921, Harbinger, NC 27941	
33. BILLING PROVIDER INFO & PH # (999) 999-9999		33. BILLING PROVIDER INFO & PH # (999) 999-9999	
SIGNED DATE		a. 5145970699 b. HSN39VRE3MZ936 c. 9046867708 d. MBX32GEX1EGM5210L	

NUCC Instruction Manual available at: www.nucc.org APPROVED OMB-0938-0999 FORM CMS-1500 (08-05)

Appendix 3 – Example of a Hand-Printed CMS-1500 Form (Synthetic Data)

1500
HEALTH INSURANCE CLAIM FORM

APPROVED BY NATIONAL UNIFORM CLAIM COMMITTEE 08/05

<input type="checkbox"/> PICA <input type="checkbox"/> PICA	
1. MEDICARE <input type="checkbox"/> MEDICAID <input type="checkbox"/> TRICARE CHAMPUS <input checked="" type="checkbox"/> CHAMPVA <input type="checkbox"/> GROUP HEALTH PLAN <input type="checkbox"/> FECA BLK (LUNG) <input type="checkbox"/> OTHER <input type="checkbox"/> <small>(Medicare #) (Medicaid #) (Sponsor's SSN) (Member ID) (SSN or ID) (SSN)</small>	
2. PATIENT'S NAME (Last Name, First Name, Middle Initial) Finch, Eiton J.	
3. PATIENT'S BIRTH DATE <input type="checkbox"/> SEX 05/06/1988 M <input checked="" type="checkbox"/> F <input type="checkbox"/>	
4. INSURED'S I.D. NUMBER (For Program in Item 1) 62SV189NOYJ332P58KUK7402UOMC	
4. INSURED'S NAME (Last Name, First Name, Middle Initial) Bray, Thor B.	
5. PATIENT'S ADDRESS (No., Street) 2447 Romig Place - Suite 771	
6. PATIENT RELATIONSHIP TO INSURED Self <input checked="" type="checkbox"/> Spouse <input type="checkbox"/> Child <input type="checkbox"/> Other <input type="checkbox"/>	
7. INSURED'S ADDRESS (No., Street) 2447 Romig Place - Suite 771	
8. PATIENT STATUS Single <input type="checkbox"/> Married <input checked="" type="checkbox"/> Other <input type="checkbox"/>	
9. OTHER INSURED'S NAME (Last Name, First Name, Middle Initial) Gardner, Yvette F.	
10. IS PATIENT'S CONDITION RELATED TO: a. EMPLOYMENT? (Current or Previous) YES <input checked="" type="checkbox"/> NO <input type="checkbox"/> b. AUTO ACCIDENT? YES <input checked="" type="checkbox"/> NO <input type="checkbox"/> PLACE (State) PA c. OTHER ACCIDENT? YES <input type="checkbox"/> NO <input checked="" type="checkbox"/>	
11. INSURED'S POLICY GROUP OR FECA NUMBER 06PP665R9CM886X9220N75KCA7HD	
12. INSURED'S DATE OF BIRTH <input type="checkbox"/> SEX 03/29/99 M <input checked="" type="checkbox"/> F <input type="checkbox"/>	
13. INSURED'S POLICY OR GROUP NUMBER 3911793X607322V42NLA53RBV5SF	
14. INSURED'S NAME OR SCHOOL NAME University of Ontario Institute of Technology	
15. INSURED'S NAME OR PROGRAM NAME Ingalls Family Care Center	
16. IS THERE ANOTHER HEALTH BENEFIT PLAN? YES <input type="checkbox"/> NO <input checked="" type="checkbox"/> If yes, return to and complete item 9 a-d	
17. PATIENT'S OR AUTHORIZED PERSON'S SIGNATURE authorizes payment of medical benefits to the undersigned physician or supplier for services described below Signed: Ross Cobb DATE: 03/07/09	
18. INSURED'S OR AUTHORIZED PERSON'S SIGNATURE authorizes payment of medical benefits to the undersigned physician or supplier for services described below Signed: Thor Bray	
19. RESERVED FOR LOCAL USE xxxxx	
20. OUTSIDE LAB? <input type="checkbox"/> YES <input checked="" type="checkbox"/> NO CHARGES 1234 23	
21. DIAGNOSIS OR NATURE OF ILLNESS OR INJURY (Relate Items 1, 2, 3 or 4 to Item 24E by Line) 1. WNS, 7745 3. xvj9, 8533 2. of74 5057 4. 2c88, 6578	
22. MEDICAID RESUBMISSION ORIGINAL REF NO. 6RP598KR70L LPV42JRU5WJ026SV98	
23. PRIOR AUTHORIZATION NUMBER Y1128BJ12KY032EL4339RB5L546WQ	
24. A. DATE(S) OF SERVICE From To B. PLACE OF SERVICE C. D. PROCEDURES, SERVICES, OR SUPPLIES (Explain Unusual Circumstances) E. DIAGNOSIS POINTER F. G. DAYS OR UNITS H. I. RENDERING PROVIDER ID # J. 1. 12/14/02 01/15/03 VI HP 55KA81 4W8C4Wwo 7745 7325d339Y NPI YJA27Y0U4SD 2. 11/02/10 12/13/09 SE EA 6E4256 2VJ42VJ4 8533 13209812dC NPI NMO24DP49QZ 3. 03/17/05 11/21/05 9B LB 27VQ11 92A8Y9a0V 5057 66952 265 NPI NAO93RRY7KY 4. 11/26/09 10/27/05 V8 D6 8Y56L0 6J0Y6J0Y 6578 56332 846 NPI KJR30c2B2KM 5. NPI 6. NPI	
25. FEDERAL TAX I.D. NUMBER <input type="checkbox"/> SSN EIN <input checked="" type="checkbox"/> 418592868023390	
26. PATIENT'S ACCOUNT NO. 0BK20VSB6MD172	
27. ACCEPT ASSIGNMENT? YES <input checked="" type="checkbox"/> NO <input type="checkbox"/>	
28. TOTAL CHARGE \$ 00 29. AMOUNT PAID \$ 00 30. BALANCE DUE \$ 00	
31. SIGNATURE OF PHYSICIAN OR SUPPLIER INCLUDING DEGREE(S) OR CREDENTIALS (I certify that the statements on the reverse apply to this bill and are made a part thereof.) Dr Ja Public 05/02/09	
32. SERVICE FACILITY LOCATION INFORMATION 800 Spruce Street, Philadelphia, PA	
33. BILLING PROVIDER INFO & PH # (999) 999-9999 1962B Haricot Drive, Apartment A2, Yellow Jacket	
SIGNED: 7234993880 DATE: DML78PBL8YPL297 6245455636 vCE37U0P7PFMZ4706Y	

NUCC Instruction Manual available at: www.nucc.org APPROVED OMB-0938-0999 FORM CMS-1500 (08-05)

Yes, this one is a bit challenging, but that was done on purpose for a client.